

Deep Recurrent Neural Network with Gated Recurrent Units

Alex Adamson
aadamson@stanford.edu

May 13, 2015

1 Network specification

This network is an extension of the network introduced in Irsoy and Cardie, 2014. The transfer functions between the units in a single layer now use gated recurrent methods.

For $i > 1$, we have

$$\vec{z}_t^{(i)} = f_2(\overrightarrow{Wz}^{(i)} \vec{h}_t^{(i-1)} + \overleftarrow{Wz}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overrightarrow{Vz}^{(i)} \vec{h}_{t-1}^{(i)}) \quad (1)$$

$$\vec{r}_t^{(i)} = f_2(\overrightarrow{Wr}^{(i)} \vec{h}_t^{(i-1)} + \overleftarrow{Wr}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overrightarrow{Vr}^{(i)} \vec{h}_{t-1}^{(i)}) \quad (2)$$

$$\widetilde{\vec{h}}_t^{(i)} = f(\overrightarrow{W}^{(i)} \vec{h}_t^{(i-1)} + \overleftarrow{W}^{(i)} \overleftarrow{h}_t^{(i-1)} + \vec{r}_t^{(i)} \circ \overrightarrow{V}^{(i)} \vec{h}_{t-1}^{(i)}) \quad (3)$$

$$\vec{h}_t^{(i)} = \vec{z}_t^{(i)} \circ \vec{h}_{t-1}^{(i)} + (1 - \vec{z}_t^{(i)}) \circ \widetilde{\vec{h}}_t^{(i)} \quad (4)$$

$$\overleftarrow{r}_t^{(i)} = f_2(\overleftarrow{Wr}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overleftarrow{W}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overleftarrow{Vr}^{(i)} \overleftarrow{h}_{t+1}^{(i)}) \quad (5)$$

$$\overleftarrow{z}_t^{(i)} = f_2(\overleftarrow{Wz}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overleftarrow{W}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overleftarrow{Vz}^{(i)} \overleftarrow{h}_{t+1}^{(i)}) \quad (6)$$

$$\widetilde{\overleftarrow{h}}_t^{(i)} = f(\overleftarrow{W}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overleftarrow{W}^{(i)} \overleftarrow{h}_t^{(i-1)} + \overleftarrow{r}_t^{(i)} \circ \overleftarrow{V}^{(i)} \overleftarrow{h}_{t+1}^{(i)}) \quad (7)$$

$$\overleftarrow{h}_t^{(i)} = \overleftarrow{z}_t^{(i)} \circ \overleftarrow{h}_{t+1}^{(i)} + (1 - \overleftarrow{z}_t^{(i)}) \circ \widetilde{\overleftarrow{h}}_t^{(i)} \quad (8)$$

and for $i = 1$ we have

$$\vec{h}_t^{(1)} = f(\overrightarrow{W}^{(1)} x_t + \overrightarrow{V}^{(1)} \vec{h}_{t-1}^{(1)} + \vec{b}^{(1)}) \quad (9)$$

$$\overleftarrow{h}_t^{(1)} = f(\overleftarrow{W}^{(1)} x_t + \overleftarrow{V}^{(1)} \overleftarrow{h}_{t+1}^{(1)} + \overleftarrow{b}^{(1)}) \quad (10)$$

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

$$y_t = g(\overrightarrow{U} \vec{h}_t^{(L)} + \overleftarrow{U} \overleftarrow{h}_t^{(L)} + c) \quad (11)$$

2 Backpropagation derivation

We again use error vector notation when deriving the gradients via backpropagation.

Let δ_t^y be the error vector propagated by the softmax unit at timestep t . Note that this is unchanged from the original network.

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 .

We first derive $\vec{\delta}_t^{(i)}$. Let f_2^* be a function such that $f_2^*(f_2(x)) = f_2'(x)$. Let f^* be the function such that $f^*(f(x)) = f'(x)$. Then,

$$\vec{\delta}_t^{(i)} = \frac{\partial \vec{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \frac{\partial J}{\partial \vec{h}_t^{(i+1)}} + \frac{\partial \overleftarrow{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \frac{\partial J}{\partial \overleftarrow{h}_t^{(i+1)}} + \frac{\partial \vec{h}_{t+1}^i}{\partial \vec{h}_t^{(i)}} \frac{\partial J}{\partial \vec{h}_{t+1}^i}$$

$$\begin{aligned} \frac{\partial \vec{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \frac{\partial J}{\partial \vec{h}_t^{(i+1)}} &= \frac{\partial \vec{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \cdot \vec{\delta}_t^{(i+1)} \\ &= \left(\frac{\partial \vec{z}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \circ \vec{h}_{t-1}^{(i+1)} + \frac{\partial (1 - \vec{z}_t^{(i+1)})}{\partial \vec{h}_t^{(i)}} \circ \overleftarrow{h}_t^{(i+1)} + (1 - \vec{z}_t^{(i+1)}) \circ \frac{\partial \overleftarrow{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \right) \cdot \vec{\delta}_t^{(i+1)} \\ &= \left(\frac{\partial \vec{z}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \circ (\vec{h}_{t-1}^{(i+1)} - \overleftarrow{h}_t^{(i+1)}) + (1 - \vec{z}_t^{(i+1)}) \circ \frac{\partial \overleftarrow{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \right) \cdot \vec{\delta}_t^{(i+1)} \\ &= ((f_2^*(\vec{z}_t^{(i+1)}) \circ (\overleftarrow{W}_{\vec{z}}^{(i+1)})^T) \circ (\vec{h}_{t-1}^{(i+1)} - \overleftarrow{h}_t^{(i+1)}) + (1 - \vec{z}_t^{(i+1)}) \circ \frac{\partial \overleftarrow{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}}) \cdot \vec{\delta}_t^{(i+1)} \\ &= ((f_2^*(\vec{z}_t^{(i+1)}) \circ (\overleftarrow{W}_{\vec{z}}^{(i+1)})^T) \circ (\vec{h}_{t-1}^{(i+1)} - \overleftarrow{h}_t^{(i+1)}) + (1 - \vec{z}_t^{(i+1)}) \circ f^*(\overleftarrow{h}_t^{(i+1)}) \circ (\overleftarrow{W}_{\vec{z}}^{(i+1)})^T) \cdot \vec{\delta}_t^{(i+1)} \end{aligned}$$

$$\begin{aligned} \frac{\partial \overleftarrow{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \frac{\partial J}{\partial \overleftarrow{h}_t^{(i+1)}} &= \frac{\partial \overleftarrow{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \cdot \overleftarrow{\delta}_t^{(i+1)} \\ &= \left(\frac{\partial \overleftarrow{z}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \circ \overleftarrow{h}_{t+1}^{(i+1)} + \frac{\partial (1 - \overleftarrow{z}_t^{(i+1)})}{\partial \vec{h}_t^{(i)}} \circ \overleftarrow{h}_t^{(i+1)} + (1 - \overleftarrow{z}_t^{(i+1)}) \circ \frac{\partial \overleftarrow{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \right) \cdot \overleftarrow{\delta}_t^{(i+1)} \\ &= \left(\frac{\partial \overleftarrow{z}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \circ (\overleftarrow{h}_{t+1}^{(i+1)} - \overleftarrow{h}_t^{(i+1)}) + (1 - \overleftarrow{z}_t^{(i+1)}) \circ \frac{\partial \overleftarrow{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}} \right) \cdot \overleftarrow{\delta}_t^{(i+1)} \\ &= ((f_2^*(\overleftarrow{z}_t^{(i+1)}) \circ (\overleftarrow{W}_{\vec{z}}^{(i+1)})^T) \circ (\overleftarrow{h}_{t+1}^{(i+1)} - \overleftarrow{h}_t^{(i+1)}) + (1 - \overleftarrow{z}_t^{(i+1)}) \circ \frac{\partial \overleftarrow{h}_t^{(i+1)}}{\partial \vec{h}_t^{(i)}}) \cdot \overleftarrow{\delta}_t^{(i+1)} \\ &= ((f_2^*(\overleftarrow{z}_t^{(i+1)}) \circ (\overleftarrow{W}_{\vec{z}}^{(i+1)})^T) \circ (\overleftarrow{h}_{t+1}^{(i+1)} - \overleftarrow{h}_t^{(i+1)}) + (1 - \overleftarrow{z}_t^{(i+1)}) \circ f^*(\overleftarrow{h}_t^{(i+1)}) \circ (\overleftarrow{W}_{\vec{z}}^{(i+1)})^T) \cdot \overleftarrow{\delta}_t^{(i+1)} \end{aligned}$$

$$\begin{aligned}
\frac{\partial \vec{h}_{t+1}^{(i)}}{\partial \vec{h}_t^{(i)}} \frac{\partial J}{\partial \vec{h}_{t+1}^{(i)}} &= \frac{\partial \vec{h}_{t+1}^{(i0)}}{\partial \vec{h}_t^{(i)}} \cdot \vec{\delta}_{t+1}^{(i)} \\
&= (\vec{z}_{t+1}^{(i)} + \frac{\partial \vec{z}_{t+1}^{(i)}}{\partial \vec{h}_t^{(i)}} \circ \vec{h}_t^{(i)} + \frac{\partial(1 - \vec{z}_{t+1}^{(i)})}{\partial \vec{h}_t^{(i)}} \circ \widetilde{\vec{h}}_{t+1}^{(i)} + (1 - \vec{z}_{t+1}^{(i)}) \circ \frac{\partial \vec{h}_{t+1}^{(i)}}{\partial \vec{h}_t^{(i)}}) \cdot \vec{\delta}_{t+1}^{(i)} \\
&= (\vec{z}_{t+1}^{(i)} + \frac{\partial \vec{z}_{t+1}^{(i)}}{\partial \vec{h}_t^{(i)}} \circ (\vec{h}_t^{(i)} - \widetilde{\vec{h}}_{t+1}^{(i)}) + (1 - \vec{z}_{t+1}^{(i)}) \circ \frac{\partial \vec{h}_{t+1}^{(i)}}{\partial \vec{h}_t^{(i)}}) \cdot \vec{\delta}_{t+1}^{(i)} \\
&= (\vec{z}_{t+1}^{(i)} + (f_2^*(\vec{z}_{t+1}^{(i)}) \circ (\vec{V}\vec{z}^{(i)})^T) \circ (\vec{h}_t^{(i)} - \widetilde{\vec{h}}_{t+1}^{(i)}) + (1 - \vec{z}_{t+1}^{(i)}) \circ \frac{\partial \vec{h}_{t+1}^{(i)}}{\partial \vec{h}_t^{(i)}}) \cdot \vec{\delta}_{t+1}^{(i)} \\
&= (\vec{z}_{t+1}^{(i)} + (f_2^*(\vec{z}_{t+1}^{(i)}) \circ (\vec{V}\vec{z}^{(i)})^T) \circ (\vec{h}_t^{(i)} - \widetilde{\vec{h}}_{t+1}^{(i)}) \\
&\quad + (1 - \vec{z}_{t+1}^{(i)}) \circ f^*(\widetilde{\vec{h}}_{t+1}^{(i)}) \circ \frac{\partial(\vec{r}_{t+1}^{(i)} \circ \vec{V}^{(i)} \vec{h}_t^{(i)})}{\partial \vec{h}_t^{(i)}}) \cdot \vec{\delta}_{t+1}^{(i)} \\
&= (\vec{z}_{t+1}^{(i)} + (f_2^*(\vec{z}_{t+1}^{(i)}) \circ (\vec{V}\vec{z}^{(i)})^T) \circ (\vec{h}_t^{(i)} - \widetilde{\vec{h}}_{t+1}^{(i)}) \\
&\quad + (1 - \vec{z}_{t+1}^{(i)}) \circ f^*(\widetilde{\vec{h}}_{t+1}^{(i)}) \circ (\frac{\partial \vec{r}_{t+1}^{(i)}}{\partial \vec{h}_t^{(i)}} + \frac{\partial \vec{V}^{(i)} \vec{h}_t^{(i)}}{\partial \vec{h}_t^{(i)}})) \cdot \vec{\delta}_{t+1}^{(i)} \\
&= (\vec{z}_{t+1}^{(i)} + (f_2^*(\vec{z}_{t+1}^{(i)}) \circ (\vec{V}\vec{z}^{(i)})^T) \circ (\vec{h}_t^{(i)} - \widetilde{\vec{h}}_{t+1}^{(i)}) \\
&\quad + (1 - \vec{z}_{t+1}^{(i)}) \circ f^*(\widetilde{\vec{h}}_{t+1}^{(i)}) \circ (f_2^*(\vec{r}_{t+1}^{(i)}) \circ (\vec{V}\vec{r}^{(i)})^T + (\vec{V}^{(i)})^T)) \cdot \vec{\delta}_{t+1}^{(i)}
\end{aligned}$$