

Gated Recurrent Unit (GRU): back-propagation derivation

Harsha Vardhan

May 20, 2022

Abstract

This document contains derivation of the gradients for a Gated Recurrent Unit (a type of RNN-block) using back-propagation (or reverse-mode differentiation). For the implementation of the neural-network see the accompanying notebooks.

1 Block Architecture

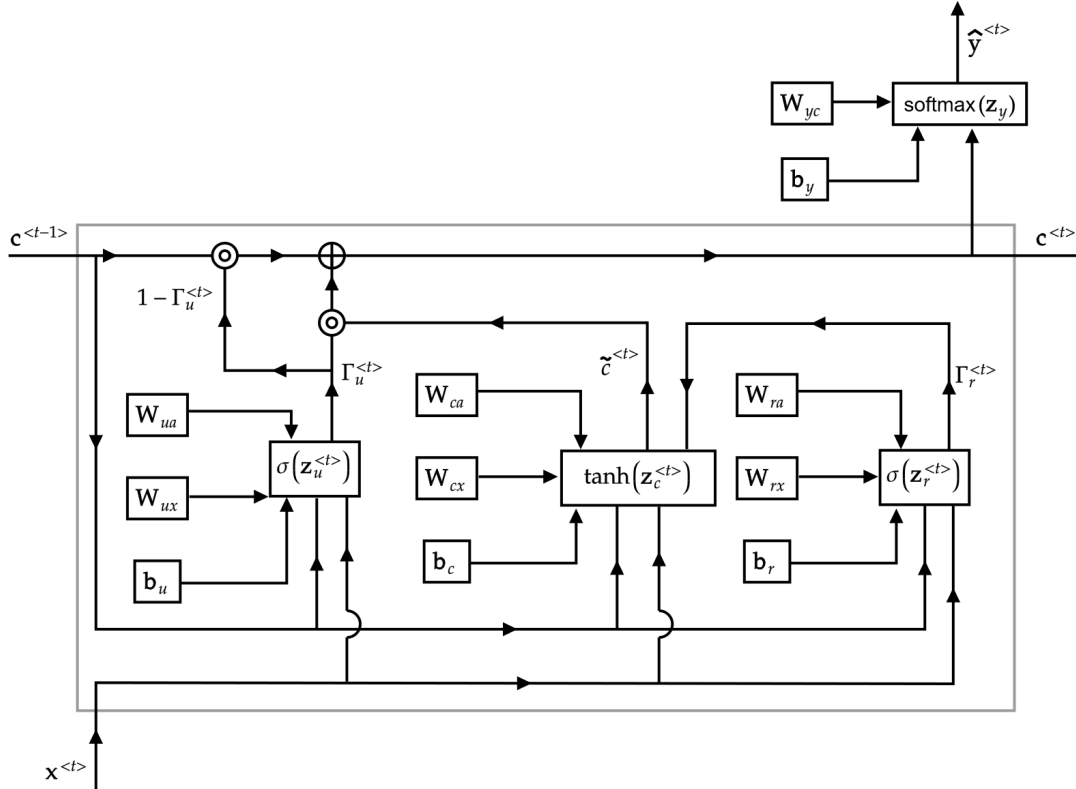


Figure 1: forward propagation diagram for a GRU-block at time-step t . The concentric circles represent a hadamard-product (i.e. $\mathbf{x} \odot \mathbf{y}$) of the input vectors.

2 Forward Propagation

Given, $\mathbf{c}^{(t-1)}$ from the $(t-1)^{th}$ time-step, the equations for forward propagation through the t^{th} time-step, are as follows: (see figure-1)

$$\mathbf{z}_u^{(t)} = \mathbf{W}_{uc}^\top \mathbf{c}^{(t-1)} + \mathbf{W}_{ux}^\top \mathbf{x}^{(t)} + \mathbf{b}_u \quad \text{eq.1}$$

$$\Gamma_u^{(t)} = \sigma(\mathbf{z}_u^{(t)}) \quad \text{eq.2}$$

$$\mathbf{z}_r^{(t)} = \mathbf{W}_{rc}^\top \mathbf{c}^{(t-1)} + \mathbf{W}_{rx}^\top \mathbf{x}^{(t)} + \mathbf{b}_r \quad \text{eq.3}$$

$$\Gamma_r^{(t)} = \sigma(\mathbf{z}_r^{(t)}) \quad \text{eq.4}$$

$$\mathbf{z}_c^{(t)} = \mathbf{W}_{cc}^\top (\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)}) + \mathbf{W}_{cx}^\top \mathbf{x}^{(t)} + \mathbf{b}_c \quad \text{eq.5}$$

$$\tilde{\mathbf{c}}^{(t)} = \tanh(\mathbf{z}_c^{(t)}) \quad \text{eq.6}$$

$$\mathbf{c}^{(t)} = \Gamma_u^{(t)} \circ \tilde{\mathbf{c}}^{(t)} + (1 - \Gamma_u^{(t)}) \circ \mathbf{c}^{(t-1)} \quad \text{eq.7}$$

$$\mathbf{z}_y^{(t)} = \mathbf{W}_{ya}^\top \mathbf{c}^{(t)} + \mathbf{b}_y \quad \text{eq.8}$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{z}_y^{(t)}) \quad \text{eq.9}$$

where,

$\mathbf{x}^{(t)}$ is a $(n_x, 1)$ -dimensional vector

$\mathbf{c}^{(t)}$ & $\mathbf{c}^{(t-1)}$ & $\tilde{\mathbf{c}}^{(t)}$ are $(n_a, 1)$ -dimensional vectors

$\mathbf{W}_{*,*}, * \in \{r, u, c\}$ are (n_a, n_a) -dimensional parameter matrices

$\mathbf{W}_{*,*}, * \in \{r, u, c\}$ are (n_x, n_a) -dimensional parameter matrices

$\mathbf{b}_*, * \in \{r, u, c\}$ are $(n_a, 1)$ -dimensional bias-vectors

$\Gamma_*^{(t)}, * \in \{r, u\}$ are $(n_a, 1)$ -dimensional vector, which represents a gate

$\mathbf{z}_*^{(t)}, * \in \{r, u, c\}$ are $(n_a, 1)$ -dimensional vector, is input to a gates' activation

\mathbf{W}_{yc} is a (n_a, n_y) -dimensional parameter matrix

\mathbf{b}_y is a $(n_y, 1)$ -dimensional bias vector

Also, the sub-scripts denote the following: $_r$ - relevance-gate, and $_u$ - update-gate.

2.1 Backward Propagation

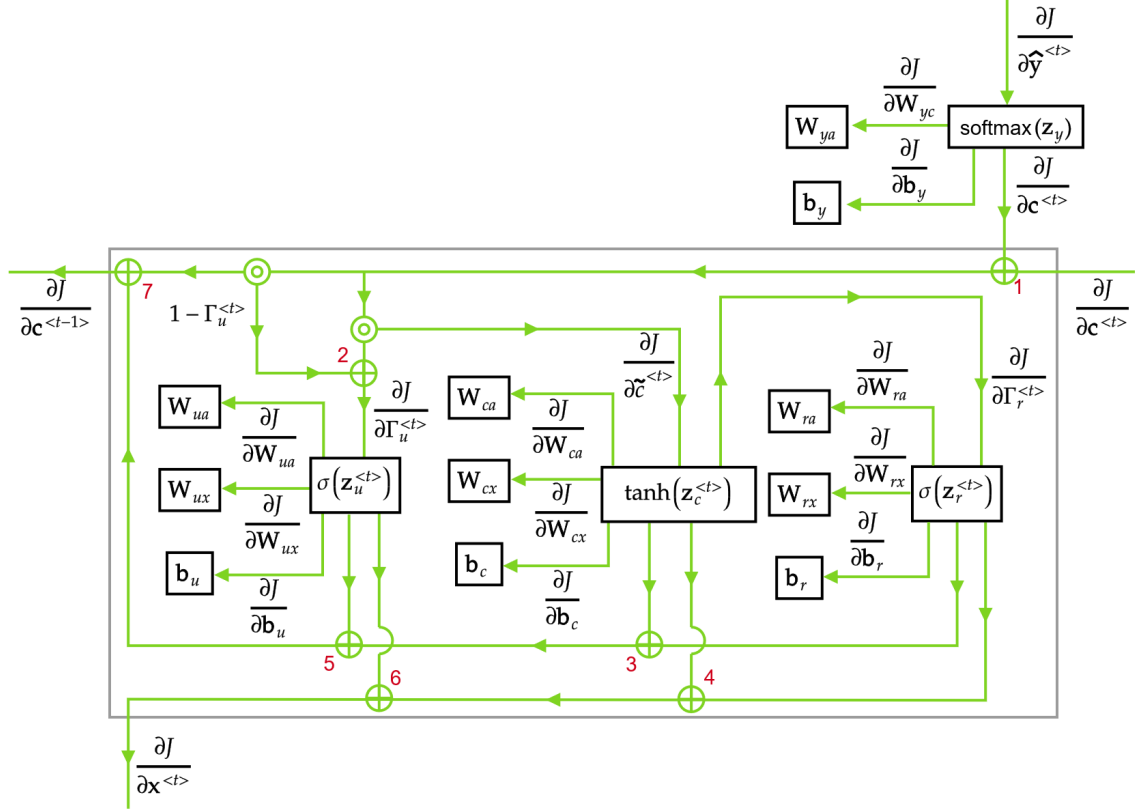


Figure 2: backward-propagation/gradient-flow diagram for a LSTM-block at time-step t . The concentric circles represent a hadamard-product (i.e. $\mathbf{x} \circ \mathbf{y}$) of the input vectors.

2.1.1 Computing $\frac{\partial J}{\partial \hat{\mathbf{y}}^{(t)}}$

Since, $\hat{\mathbf{y}}^{(t)}$ is computed using the $\text{softmax}()$ activation-function, the loss J is computed using the cross-entropy loss. Also, let $\mathbf{y}^{(t)}$ be the output-label corresponding to the t^{th} time-step. Then,

$$\begin{aligned} \frac{\partial J}{\partial \hat{\mathbf{y}}^{(t)}} &= \begin{bmatrix} \frac{\partial J}{\partial \hat{y}_1^{(t)}} & \frac{\partial J}{\partial \hat{y}_2^{(t)}} & \cdots & \frac{\partial J}{\partial \hat{y}_{n_y}^{(t)}} \end{bmatrix} \\ &= \begin{bmatrix} \frac{y_1^{(t)}}{\hat{y}_1^{(t)}} & \frac{y_2^{(t)}}{\hat{y}_2^{(t)}} & \cdots & \frac{y_{n_y}^{(t)}}{\hat{y}_{n_y}^{(t)}} \end{bmatrix} \end{aligned}$$

2.1.2 Computing $\frac{\partial J}{\partial \mathbf{W}_{yc}}$ and $\frac{\partial J}{\partial \mathbf{b}_y}$

From *eq.9* (see section-2), we have

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{z}_y^{(t)}} &= \frac{\partial J}{\partial \hat{\mathbf{y}}^{(t)}} \frac{\partial \hat{\mathbf{y}}^{(t)}}{\partial \mathbf{z}_y^{(t)}} \\ &= \frac{\partial J}{\partial \hat{\mathbf{y}}^{(t)}} \left(\text{diag}(\hat{\mathbf{y}}^{(t)}) - \hat{\mathbf{y}}^{(t)} (\hat{\mathbf{y}}^{(t)})^T \right) \end{aligned}$$

Note: for more on the derivation of the derivative of a `softmax()` function, see `..\notes\softmax-function.ipynb`

From *eq.8* (see section-2), we have

$$\text{tr}(J) = \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_y^{(t)}} d\mathbf{z}_y^{(t)} \right)$$

where, $d\mathbf{z}_y^{(t)}$ can be expanded as follows:

$$\begin{aligned} d\mathbf{z}_y^{(t)} &= d(\mathbf{W}_{yc}^\top \mathbf{c}^{(t)} + \mathbf{b}_y) \\ &= d(\mathbf{W}_{yc}^\top) \mathbf{c}^{(t)} + \mathbf{W}_{yc}^\top d\mathbf{c}^{(t)} + d\mathbf{b}_y \end{aligned} \quad \text{eq.9-1}$$

when differentiating w.r.t. \mathbf{W}_{yc} , we have $d\mathbf{c}^{(t)} = 0$, and $d\mathbf{b}_y = 0$. So,

$$\begin{aligned} dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_y^{(t)}} d\mathbf{W}_{yc}^\top \mathbf{c}^{(t)} \right) \\ &= \text{tr} \left((\mathbf{c}^{(t)})^\top d\mathbf{W}_{yc} \left(\frac{\partial J}{\partial \mathbf{z}_y^{(t)}} \right)^\top \right) \\ &= \text{tr} \left(\left(\frac{\partial J}{\partial \mathbf{z}_y^{(t)}} \right)^\top (\mathbf{c}^{(t)})^\top d\mathbf{W}_{yc} \right) \\ &\Rightarrow \frac{\partial J}{\partial \mathbf{W}_{yc}} = \left(\frac{\partial J}{\partial \mathbf{z}_y^{(t)}} \right)^\top (\mathbf{c}^{(t)})^\top \end{aligned}$$

when differentiating w.r.t. \mathbf{b}_y , we have $d\mathbf{c}^{(t)} = 0$, and $d\mathbf{W}_{yc} = 0$. So,

$$\begin{aligned} dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_y^{(t)}} d\mathbf{b}_y^\top \right) \\ &\Rightarrow \frac{\partial J}{\partial \mathbf{b}_y} = \frac{\partial J}{\partial \mathbf{z}_y^{(t)}} \end{aligned}$$

2.1.3 Computing $\frac{\partial J}{\partial \mathbf{c}^{(t)}}$

This derivative has two components

Comp-1 flows-in from the $(t+1)^{th}$ time-step, and

Comp-2 flows-in as the derivative from $\hat{\mathbf{y}}^{(t)}$. This derivative is computed using *eq.9-1* as follows,

$$\begin{aligned} dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_y^{(t)}} d\mathbf{z}_y^{(t)} \right) \\ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_y^{(t)}} \mathbf{W}_{yc}^\top d\mathbf{c}^{(t)} \right) \\ &\Rightarrow \frac{\partial J}{\partial \mathbf{c}^{(t)}} = \frac{\partial J}{\partial \mathbf{z}_y^{(t)}} \mathbf{W}_{yc}^\top \end{aligned}$$

These two components are added to compute the true derivative (see the \oplus labeled as 1 in the figure-2), i.e.

$$\frac{\partial J}{\partial \mathbf{c}^{(t)}} = \left. \frac{\partial J}{\partial \mathbf{c}^{(t)}} \right|_{\text{Comp-1}} + \left. \frac{\partial J}{\partial \mathbf{c}^{(t)}} \right|_{\text{Comp-2}}$$

2.1.4 Computing $\frac{\partial J}{\partial \tilde{\mathbf{c}}^{(t)}}$ and $\frac{\partial J}{\partial \Gamma_u^{(t)}}$

From eq.7 (see section-2), we have

$$\begin{aligned} \frac{\partial J}{\partial \tilde{\mathbf{c}}^{(t)}} &= \frac{\partial J}{\partial \mathbf{c}^{(t)}} \frac{\partial \mathbf{c}^{(t)}}{\partial \tilde{\mathbf{c}}^{(t)}} \\ &= \frac{\partial J}{\partial \mathbf{c}^{(t)}} \text{diag}(\Gamma_u^{(t)}) \\ &= \frac{\partial J}{\partial \mathbf{c}^{(t)}} \circ (\Gamma_u^{(t)})^\top \\ \frac{\partial J}{\partial \Gamma_u^{(t)}} &= \frac{\partial J}{\partial \mathbf{c}^{(t)}} \frac{\partial \mathbf{c}^{(t)}}{\partial \Gamma_u^{(t)}} \\ &= \frac{\partial J}{\partial \mathbf{c}^{(t)}} \text{diag}(\tilde{\mathbf{c}}^{(t)} - \mathbf{c}^{(t-1)}) \\ &= \frac{\partial J}{\partial \mathbf{c}^{(t)}} \circ (\tilde{\mathbf{c}}^{(t)} - \mathbf{c}^{(t-1)})^\top \end{aligned}$$

2.1.5 Computing $\frac{\partial J}{\partial \mathbf{W}_{cc}}$, $\frac{\partial J}{\partial \mathbf{W}_{cx}}$, and $\frac{\partial J}{\partial \mathbf{b}_c}$

From eq.6 (see section-2), we have

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{z}_c^{(t)}} &= \frac{\partial J}{\partial \tilde{\mathbf{c}}^{(t)}} \frac{\partial \tilde{\mathbf{c}}^{(t)}}{\partial \mathbf{z}_c^{(t)}} \\ &= \frac{\partial J}{\partial \tilde{\mathbf{c}}^{(t)}} \circ [\mathbf{1}_{(n_a,1)} - \tanh^2(\mathbf{z}_c^{(t)})]^\top \end{aligned}$$

From eq.5 (see section-2), using the trace-method, we have

$$\text{tr}(J) = \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} d\mathbf{z}_c^{(t)} \right)$$

where, $d\mathbf{z}_c^{(t)}$ can be expanded as follows:

$$\begin{aligned} d\mathbf{z}_c^{(t)} &= d(\mathbf{W}_{cc}^\top (\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)}) + \mathbf{W}_{cx}^\top \mathbf{x}^{(t)} + \mathbf{b}_c) \\ &= d\mathbf{W}_{cc}^\top (\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)}) + \mathbf{W}_{cc}^\top d(\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)}) + d\mathbf{W}_{cx}^\top \mathbf{x}^{(t)} \\ &\quad + \mathbf{W}_{cx}^\top d\mathbf{x}^{(t)} + d\mathbf{b}_c \quad \text{eq.9-2} \end{aligned}$$

when differentiating w.r.t. \mathbf{W}_{cc} , we have $d\mathbf{W}_{cx} = 0$, $d(\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)}) = 0$, $d\mathbf{b}_c = 0$, and $d\mathbf{x}^{(t)} = 0$. So,

$$\begin{aligned}
dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} d\mathbf{W}_{cc}^\top (\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)}) \right) \\
&= \text{tr} \left((\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)})^\top d\mathbf{W}_{cc} \left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \right)^\top \right) \\
&= \text{tr} \left(\left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \right)^\top (\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)})^\top d\mathbf{W}_{cc} \right) \\
&\Rightarrow \frac{\partial J}{\partial \mathbf{W}_{cc}} = \left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \right)^\top (\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)})^\top
\end{aligned}$$

when differentiating w.r.t. \mathbf{W}_{cx} , we have $d\mathbf{W}_{cc} = 0$, $d(\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)}) = 0$, $d\mathbf{b}_c = 0$, and $d\mathbf{x}^{(t)} = 0$. So,

$$\begin{aligned}
dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} d\mathbf{W}_{cx}^\top \mathbf{x}^{(t)} \right) \\
&= \text{tr} \left((\mathbf{x}^{(t)})^\top d\mathbf{W}_{cx} \left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \right)^\top \right) \\
&= \text{tr} \left(\left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \right)^\top (\mathbf{x}^{(t)})^\top d\mathbf{W}_{cx} \right) \\
&\Rightarrow \frac{\partial J}{\partial \mathbf{W}_{cx}} = \left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \right)^\top (\mathbf{x}^{(t)})^\top
\end{aligned}$$

when differentiating w.r.t. \mathbf{b}_c , we have $d\mathbf{W}_{cc} = 0$, $d\mathbf{c}^{(t-1)} = 0$, $d\mathbf{W}_{cx} = 0$, and $d\mathbf{x}^{(t)} = 0$. So,

$$\begin{aligned}
dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} d\mathbf{b}_c \right) \\
&\Rightarrow \frac{\partial J}{\partial \mathbf{b}_c} = \frac{\partial J}{\partial \mathbf{z}_c^{(t)}}
\end{aligned}$$

2.1.6 Computing $\frac{\partial J}{\partial \mathbf{W}_{uc}}$, $\frac{\partial J}{\partial \mathbf{W}_{ux}}$, and $\frac{\partial J}{\partial \mathbf{b}_u}$

From *eq.2* (see section-2), we have

$$\begin{aligned}
\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} &= \frac{\partial J}{\partial \Gamma_u^{(t)}} \frac{\partial \Gamma_u^{(t)}}{\partial \mathbf{z}_u^{(t)}} \\
&= \frac{\partial J}{\partial \Gamma_u^{(t)}} \circ (\Gamma_u^{(t)} \circ (\mathbf{1}_{(n_c, 1)} - \Gamma_u^{(t)}))^\top
\end{aligned}$$

From *eq.1* (see section-2), using the trace-method, we have

$$\text{tr}(J) = \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} d\mathbf{z}_u^{(t)} \right)$$

where, $d\mathbf{z}_u^{(t)}$ can be expanded as follows:

$$\begin{aligned}
d\mathbf{z}_u^{(t)} &= d(\mathbf{W}_{uc}^\top \mathbf{c}^{(t-1)} + \mathbf{W}_{ux}^\top \mathbf{x}^{(t)} + \mathbf{b}_u) \\
&= d\mathbf{W}_{uc}^\top \mathbf{c}^{(t-1)} + \mathbf{W}_{uc}^\top d\mathbf{c}^{(t-1)} + d\mathbf{W}_{ux}^\top \mathbf{x}^{(t)} + \mathbf{W}_{ux}^\top d\mathbf{x}^{(t)} + d\mathbf{b}_u
\end{aligned} \tag{eq.9-2}$$

when differentiating w.r.t. \mathbf{W}_{uc} , we have $d\mathbf{W}_{ux} = 0$, $d\mathbf{c}^{(t-1)} = 0$, $d\mathbf{b}_u = 0$, and $d\mathbf{x}^{(t)} = 0$. So,

$$\begin{aligned} dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} d\mathbf{W}_{uc}^\top \mathbf{c}^{(t-1)} \right) \\ &= \text{tr} \left((\mathbf{c}^{(t-1)})^\top d\mathbf{W}_{uc} \left(\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} \right)^\top \right) \\ &= \text{tr} \left(\left(\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} \right)^\top (\mathbf{c}^{(t-1)})^\top d\mathbf{W}_{uc} \right) \\ &\Rightarrow \frac{\partial J}{\partial \mathbf{W}_{uc}} = \left(\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} \right)^\top (\mathbf{c}^{(t-1)})^\top \end{aligned}$$

when differentiating w.r.t. \mathbf{W}_{ux} , we have $d\mathbf{W}_{uc} = 0$, $d\mathbf{c}^{(t-1)} = 0$, $d\mathbf{b}_u = 0$, and $d\mathbf{x}^{(t)} = 0$. So,

$$\begin{aligned} dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} d\mathbf{W}_{ux}^\top \mathbf{x}^{(t)} \right) \\ &= \text{tr} \left((\mathbf{x}^{(t)})^\top d\mathbf{W}_{ux} \left(\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} \right)^\top \right) \\ &= \text{tr} \left(\left(\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} \right)^\top (\mathbf{x}^{(t)})^\top d\mathbf{W}_{ux} \right) \\ &\Rightarrow \frac{\partial J}{\partial \mathbf{W}_{ux}} = \left(\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} \right)^\top (\mathbf{x}^{(t)})^\top \end{aligned}$$

when differentiating w.r.t. \mathbf{b}_u , we have $d\mathbf{W}_{uc} = 0$, $d\mathbf{c}^{(t-1)} = 0$, $d\mathbf{W}_{ux} = 0$, and $d\mathbf{x}^{(t)} = 0$. So,

$$\begin{aligned} dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_u^{(t)}} d\mathbf{b}_u \right) \\ &\Rightarrow \frac{\partial J}{\partial \mathbf{b}_u} = \frac{\partial J}{\partial \mathbf{z}_u^{(t)}} \end{aligned}$$

2.1.7 Computing $\frac{\partial J}{\partial \Gamma_r^{(t)}}$

From the results in section-2.1.5 and, from *eq.9-2*, we have

$$\begin{aligned} \frac{\partial J}{\partial \Gamma_r^{(t)}} &= \frac{\partial J}{\partial \tilde{\mathbf{c}}^{(t)}} \frac{\partial \tilde{\mathbf{c}}^{(t)}}{\partial \Gamma_r^{(t)}} \\ &= \frac{\partial J}{\partial \tilde{\mathbf{c}}^{(t)}} \frac{\partial \tilde{\mathbf{c}}^{(t)}}{\partial (\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)})} \frac{\partial (\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)})}{\partial \Gamma_r^{(t)}} \\ &= \frac{\partial J}{\partial \tilde{\mathbf{c}}^{(t)}} \mathbf{W}_{cc}^\top \text{diag}(\mathbf{c}^{(t-1)}) \\ &= \left(\frac{\partial J}{\partial \tilde{\mathbf{c}}^{(t)}} \mathbf{W}_{cc}^\top \right) \circ (\mathbf{c}^{(t-1)})^\top \end{aligned}$$

Note: for more information on the derivative of a Hadamard product, see the Appendix-A in `.\backprop-lstm.pdf`.

2.1.8 Computing $\frac{\partial J}{\partial \mathbf{W}_{rc}}$, $\frac{\partial J}{\partial \mathbf{W}_{rx}}$, and $\frac{\partial J}{\partial \mathbf{b}_r}$

From *eq.4* (see section-2), we have

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} &= \frac{\partial J}{\partial \Gamma_r^{(t)}} \frac{\partial \Gamma_r^{(t)}}{\partial \mathbf{z}_r^{(t)}} \\ &= \frac{\partial J}{\partial \Gamma_r^{(t)}} \circ (\Gamma_r^{(t)} \circ (\mathbf{1}_{(n_c,1)} - \Gamma_r^{(t)}))^\top\end{aligned}$$

From *eq.3* (see section-2), using the trace-method, we have

$$\text{tr}(J) = \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} d\mathbf{z}_r^{(t)} \right)$$

where, $d\mathbf{z}_r^{(t)}$ can be expanded as follows:

$$\begin{aligned}d\mathbf{z}_r^{(t)} &= d(\mathbf{W}_{rc}^\top \mathbf{c}^{(t-1)} + \mathbf{W}_{rx}^\top \mathbf{x}^{(t)} + \mathbf{b}_r) \\ &= d\mathbf{W}_{rc}^\top \mathbf{c}^{(t-1)} + \mathbf{W}_{rc}^\top d\mathbf{c}^{(t-1)} + d\mathbf{W}_{rx}^\top \mathbf{x}^{(t)} + \mathbf{W}_{rx}^\top d\mathbf{x}^{(t)} + d\mathbf{b}_r\end{aligned}\quad \text{eq.9-3}$$

when differentiating w.r.t. \mathbf{W}_{rc} , we have $d\mathbf{W}_{rx} = 0$, $d\mathbf{c}^{(t-1)} = 0$, $d\mathbf{b}_r = 0$, and $d\mathbf{x}^{(t)} = 0$. So,

$$\begin{aligned}dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} d\mathbf{W}_{rc}^\top \mathbf{c}^{(t-1)} \right) \\ &= \text{tr} \left((\mathbf{c}^{(t-1)})^\top d\mathbf{W}_{rc} \left(\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \right)^\top \right) \\ &= \text{tr} \left(\left(\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \right)^\top (\mathbf{c}^{(t-1)})^\top d\mathbf{W}_{rc} \right) \\ \implies \frac{\partial J}{\partial \mathbf{W}_{rc}} &= \left(\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \right)^\top (\mathbf{c}^{(t-1)})^\top\end{aligned}$$

when differentiating w.r.t. \mathbf{W}_{rx} , we have $d\mathbf{W}_{rc} = 0$, $d\mathbf{c}^{(t-1)} = 0$, $d\mathbf{b}_r = 0$, and $d\mathbf{x}^{(t)} = 0$. So,

$$\begin{aligned}dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} d\mathbf{W}_{rx}^\top \mathbf{x}^{(t)} \right) \\ &= \text{tr} \left((\mathbf{x}^{(t)})^\top d\mathbf{W}_{rx} \left(\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \right)^\top \right) \\ &= \text{tr} \left(\left(\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \right)^\top (\mathbf{x}^{(t)})^\top d\mathbf{W}_{rx} \right) \\ \implies \frac{\partial J}{\partial \mathbf{W}_{rx}} &= \left(\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \right)^\top (\mathbf{x}^{(t)})^\top\end{aligned}$$

when differentiating w.r.t. \mathbf{b}_r , we have $d\mathbf{W}_{rc} = 0$, $d\mathbf{c}^{(t-1)} = 0$, $d\mathbf{W}_{rx} = 0$, and $d\mathbf{x}^{(t)} = 0$. So,

$$\begin{aligned}dJ &= \text{tr} \left(\frac{\partial J}{\partial \mathbf{z}_r^{(t)}} d\mathbf{b}_r \right) \\ \implies \frac{\partial J}{\partial \mathbf{b}_r} &= \frac{\partial J}{\partial \mathbf{z}_r^{(t)}}\end{aligned}$$

2.1.9 Computing $\frac{\partial J}{\partial \mathbf{c}^{(t-1)}}$, and $\frac{\partial J}{\partial \mathbf{x}^{(t)}}$

The derivative $\frac{\partial J}{\partial \mathbf{c}^{(t-1)}}$ has four components, as follows

Comp-1 flows-in as part of the derivative $\frac{\partial J}{\partial \mathbf{c}^{(t)}}$. This derivative is computed as follows

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{c}^{(t-1)}} &= \frac{\partial J}{\partial \mathbf{c}^{(t)}} \frac{\partial \mathbf{c}^{(t)}}{\partial \mathbf{c}^{(t-1)}} \\ &= \frac{\partial J}{\partial \mathbf{c}^{(t)}} \text{diag}(\mathbf{1}_{(n_c,1)} - \Gamma_u^{(t)}) \\ &= \frac{\partial J}{\partial \mathbf{c}^{(t)}} \circ (\mathbf{1}_{(n_c,1)} - \Gamma_u^{(t)})^\top\end{aligned}$$

Comp-2 flows-in as derivative from $\Gamma_r^{(t)}$ (see *eq.9-3*), and is computed as follows

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{c}^{(t-1)}} &= \frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \frac{\partial \mathbf{z}_r^{(t)}}{\partial \mathbf{c}^{(t-1)}} \\ &= \frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \mathbf{W}_{rc}^\top\end{aligned}$$

Comp-3 flows-in as derivative from $\Gamma_u^{(t)}$ (see *eq.9-2*), and is computed as follows

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{c}^{(t-1)}} &= \frac{\partial J}{\partial \mathbf{z}_u^{(t)}} \frac{\partial \mathbf{z}_u^{(t)}}{\partial \mathbf{c}^{(t-1)}} \\ &= \frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \mathbf{W}_{uc}^\top\end{aligned}$$

Comp-4 flows-in as derivative from $\tilde{\mathbf{c}}^{(t)}$ (see *eq.9-2*), and is computed as follows

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{c}^{(t-1)}} &= \frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \frac{\partial \mathbf{z}_c^{(t)}}{\partial \mathbf{c}^{(t-1)}} \\ &= \frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \frac{\partial \mathbf{z}_c^{(t)}}{\partial (\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)})} \frac{\partial (\Gamma_r^{(t)} \circ \mathbf{c}^{(t-1)})}{\partial \mathbf{c}^{(t-1)}} \\ &= \frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \mathbf{W}_{cc}^\top \text{diag}(\Gamma_r^{(t)}) \\ &= \left(\frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \mathbf{W}_{cc}^\top \right) \circ (\Gamma_r^{(t)})^\top\end{aligned}$$

These four components are then added to compute the true derivative, i.e.

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{c}^{(t-1)}} &= \left. \frac{\partial J}{\partial \mathbf{c}^{(t-1)}} \right|_{\text{Comp-1}} + \left. \frac{\partial J}{\partial \mathbf{c}^{(t-1)}} \right|_{\text{Comp-2}} + \left. \frac{\partial J}{\partial \mathbf{c}^{(t-1)}} \right|_{\text{Comp-3}} + \left. \frac{\partial J}{\partial \mathbf{c}^{(t-1)}} \right|_{\text{Comp-4}} \\ &\quad (\text{see } \oplus \text{ labeled as } 3, 5, \text{ and } 7 \text{ in figure-2})\end{aligned}$$

The derivative $\frac{\partial J}{\partial \mathbf{c}^{(t-1)}}$ has three components, as follows

Comp-1 flows-in as derivative from $\Gamma_r^{(t)}$ (see *eq.9-3*), and is computed as follows

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{x}^{(t)}} &= \frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \frac{\partial \mathbf{z}_r^{(t)}}{\partial \mathbf{x}^{(t)}} \\ &= \frac{\partial J}{\partial \mathbf{z}_r^{(t)}} \mathbf{W}_{rx}^\top\end{aligned}$$

Comp-2 flows-in as derivative from $\Gamma_u^{(t)}$ (see *eq.9-2*), and is computed as follows

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{x}^{(t)}} &= \frac{\partial J}{\partial \mathbf{z}_u^{(t)}} \frac{\partial \mathbf{z}_u^{(t)}}{\partial \mathbf{x}^{(t)}} \\ &= \frac{\partial J}{\partial \mathbf{z}_u^{(t)}} \mathbf{W}_{ux}^\top\end{aligned}$$

Comp-3 flows-in as derivative from $\tilde{\mathbf{c}}^{(t)}$ (see *eq.9-2*), and is computed as follows

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{x}^{(t)}} &= \frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \frac{\partial \mathbf{z}_c^{(t)}}{\partial \mathbf{x}^{(t)}} \\ &= \frac{\partial J}{\partial \mathbf{z}_c^{(t)}} \mathbf{W}_{cx}^\top\end{aligned}$$

These three components are then added to compute the true derivative, i.e.

$$\begin{aligned}\frac{\partial J}{\partial \mathbf{x}^{(t)}} &= \left. \frac{\partial J}{\partial \mathbf{x}^{(t)}} \right|_{\text{Comp-1}} + \left. \frac{\partial J}{\partial \mathbf{x}^{(t)}} \right|_{\text{Comp-2}} + \left. \frac{\partial J}{\partial \mathbf{x}^{(t)}} \right|_{\text{Comp-3}} \\ &\quad (\text{see } \oplus \text{ labeled as } 4, \text{ and } 6 \text{ in figure-2})\end{aligned}$$

3 Gradient or Jacobian?

In the above derivations, we have used the numerator layout while performing matrix-derivatives. One of the consequences of this decision is that the derivatives that we have computed are in-fact jacobians and not gradients. Fortunately, gradients are just transpose of jacobians.