

Error comes from two parts, one from the next unit's back propagation  $d\bar{a}^{(t)}$  and  $d\bar{c}^{(t)}$ , and another is from the loss function of present unit  $d\hat{a}^{(t)}$  and  $d\hat{c}^{(t)}$ .

Thus:

$$da^{(t)} = d\bar{a}^{(t)} + d\hat{a}^{(t)}$$

$$dc^{(t)} = d\bar{c}^{(t)} + d\hat{c}^{(t)}$$

$$d\hat{c}^{(t)} = da^{(t)} * \Gamma_o^{(t)} * (1 - (\tanh(c^{(t)}))^2)$$

Thus:

$$dc^{(t)} = d\bar{c}^{(t)} + da^{(t)} * \Gamma_o^{(t)} * (1 - \tanh(c^{(t)})^2)$$

$$d\Gamma_o^{(t)} = da^{(t)} * \tanh(c^{(t)})$$

$$d\bar{c}^{(t)} = dc^{(t)} * \Gamma_u^{(t)}$$

$$d\Gamma_f^{(t)} = dc^{(t)} * c^{(t-1)}$$

$$d\Gamma_u^{(t)} = dc^{(t)} * \tilde{c}^{(t)}$$

Computing the gradients of parameters:

$$\begin{aligned} dw_o &= d\Gamma_o^{(t)} * \Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}) \cdot \begin{pmatrix} a^{(t-1)} \\ x^{(t)} \end{pmatrix}^T \\ &= da^{(t)} * \tanh(c^{(t)}) * \Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}) \cdot \begin{pmatrix} a^{(t-1)} \\ x^{(t)} \end{pmatrix}^T \end{aligned}$$

$$\begin{aligned} db_o &= \sum_{axis=1} d\Gamma_o^{(t)} * \Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}) \\ &= \sum_{axis=1} da^{(t)} * \tanh(c^{(t)}) * \Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}) \end{aligned}$$

$$\begin{aligned} dw_f &= d\Gamma_f^{(t)} * \Gamma_f^{(t)} * (1 - \Gamma_f^{(t)}) \cdot \begin{pmatrix} a^{(t-1)} \\ x^{(t)} \end{pmatrix}^T \\ &= dc^{(t)} * c^{(t-1)} * \Gamma_f^{(t)} * (1 - \Gamma_f^{(t)}) \cdot \begin{pmatrix} a^{(t-1)} \\ x^{(t)} \end{pmatrix}^T \end{aligned}$$

$$\begin{aligned} db_f &= \sum_{axis=1} d\Gamma_f^{(t)} * \Gamma_f^{(t)} * (1 - \Gamma_f^{(t)}) \\ &= \sum_{axis=1} dc^{(t)} * c^{(t-1)} * \Gamma_f^{(t)} * (1 - \Gamma_f^{(t)}) \end{aligned}$$

$$\begin{aligned} dw_u &= d\Gamma_u^{(t)} * \Gamma_u^{(t)} * (1 - \Gamma_u^{(t)}) \cdot \begin{pmatrix} a^{(t-1)} \\ x^{(t)} \end{pmatrix}^T \\ &= dc^{(t)} * \tilde{c}^{(t)} * \Gamma_u^{(t)} * (1 - \Gamma_u^{(t)}) \cdot \begin{pmatrix} a^{(t-1)} \\ x^{(t)} \end{pmatrix}^T \end{aligned}$$

$$\begin{aligned}
db_u &= \sum_{axis=1} d\Gamma_u^{(t)} * \Gamma_u^{(t)} * (1 - \Gamma_u^{(t)}) \\
&= \sum_{axis=1} dc^{(t)} * \tilde{c}^{(t)} * \Gamma_u^{(t)} * (1 - \Gamma_u^{(t)}) \\
dw_c &= d\tilde{c}^{(t)} * (1 - (\tilde{c}^{(t)})^2) \cdot \begin{pmatrix} a^{(t-1)} \\ x^{(t)} \end{pmatrix}^T \\
&= dc^{(t)} * \Gamma_u^{(t)} * (1 - (\tilde{c}^{(t)})^2) \cdot \begin{pmatrix} a^{(t-1)} \\ x^{(t)} \end{pmatrix}^T \\
db_c &= \sum_{axis=1} d\tilde{c}^{(t)} * (1 - (\tilde{c}^{(t)})^2) \\
&= \sum_{axis=1} dc^{(t)} * \Gamma_u^{(t)} * (1 - (\tilde{c}^{(t)})^2)
\end{aligned}$$

Computing the gradients of hidden state, previous memory state and input:

$$\begin{aligned}
dc^{(t-1)} &= dc^{(t)} * \Gamma_f^{(t)} \\
da^{(t-1)} &= d\Gamma_o^{(t)} * \Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}) * \hat{w}_o^T + d\Gamma_f^{(t)} * \Gamma_f^{(t)} * (1 - \Gamma_f^{(t)}) * \hat{w}_f^T \\
&\quad + d\Gamma_u^{(t)} * \Gamma_u^{(t)} * (1 - \Gamma_u^{(t)}) * \hat{w}_u^T + d\tilde{c}^{(t)} * (1 - (\tilde{c}^{(t)})^2) * \hat{w}_c^T \\
dx^{(t)} &= d\Gamma_o^{(t)} * \Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}) * \tilde{w}_o^T + d\Gamma_f^{(t)} * \Gamma_f^{(t)} * (1 - \Gamma_f^{(t)}) * \tilde{w}_f^T \\
&\quad + d\Gamma_u^{(t)} * \Gamma_u^{(t)} * (1 - \Gamma_u^{(t)}) * \tilde{w}_u^T + d\tilde{c}^{(t)} * (1 - (\tilde{c}^{(t)})^2) * \tilde{w}_c^T
\end{aligned}$$

Where  $\hat{w}^T$  denote  $w[:, : n_a]$ ,  $\tilde{w}^T$  denote  $w[:, n_a : n_a + n_x]$