

Recurrent Attention Unit with Step-down Optimized by NADAM

Enrique Boswell Nueve IV

March 26, 2019

In this document, the author will be providing full deviations for a Recurrent Attention Unit[2]. This unit has been modified with a Step-down Neural network to allow for differences in dimensionality between the inputs and targets. The parameters for the RAU are being updated with the optimization technique NADAM (Nesterov-accelerated adaptive moment estimation)[1].

References

- [1] Dozat, Timothy. "Incorporating nesterov momentum into adam." (2016).
- [2] Zhong, Guoqiang, Guohua Yue, and Xiao Ling. "Recurrent Attention Unit." arXiv preprint arXiv:1810.12754 (2018).

Contents

1	Forward Pass at Time t	3
1.1	Update Gate	3
1.2	Reset Gate	3
1.3	Candidate Gate	3
1.4	Attention Gate	3
1.5	Probabilistic Vector	3
1.6	Learning Function: General	4
1.7	State Vector	4
1.8	Step Down Network	4
1.9	Error Calculation	4
2	Backpropagation of Weights at Time t	4
2.1	Update Gate Present Weights	4
2.2	Update Gate Memory Weights	4
2.3	Candidate Gate Present Weights	5
2.4	Candidate Gate Memory Weights	5
2.5	Reset Gate Present Weights	5
2.6	Reset Gate Memory Weights	5
2.7	Attention Gate Weight	5
2.8	Learning Function Weight	6
2.9	Step-down Network Weight	6
3	Backpropagation of State at time t	6
3.1	Past State of Update Gate	6
3.2	Past State of Candidate Gate	6
3.3	Past State of Reset Gate	6
3.4	Past State of Attention Gate	7
3.5	Past State of State Vector	7
3.6	Backpropagation State Value	7
4	Parameter Update	7
4.1	Update Update Gate	7
4.2	Update Reset Gate	7
4.3	Update Candidate Gate	7
4.4	Update Attention Weight	8
4.5	Update Learning Function Weight	8
4.6	Optimization Technique: NADAM	8
5	Activation Functions	8
5.1	Sigmoid	8
5.2	Hyperbolic Tangent	8
5.3	Softmax	8

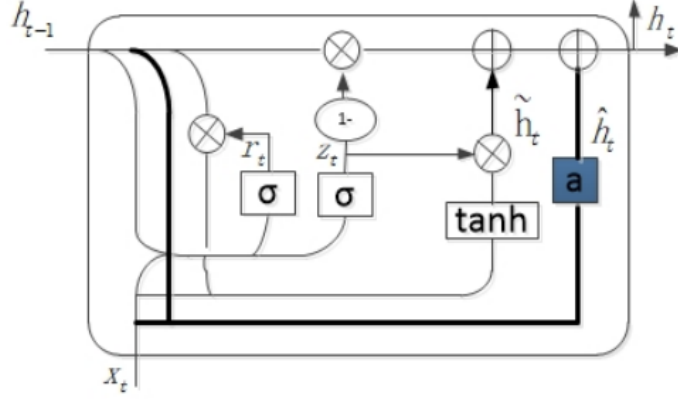


Figure 1: Recurrent Attention Unit Diagram[2].

1 Forward Pass at Time t

1.1 Update Gate

$$z_t = \sigma(W_z x_t + U_z h_{t-1})$$

1.2 Reset Gate

$$r_t = \sigma(W_r x_t + U_r h_{t-1})$$

1.3 Candidate Gate

$$\tilde{h}_t = \tanh(W_c x_t + (U_c h_{t-1} \odot r_t))$$

1.4 Attention Gate

$$\hat{h}_t = \tanh(W_a \gamma_t)$$

1.5 Probabilistic Vector

$$\gamma_t = \text{softmax}(\alpha_t)$$

1.6 Learning Function: General

$$\alpha_t = W_{\alpha_t} h_{t-1}^T x_t$$

1.7 State Vector

$$h_t = (1 - z_t)h_{t-1} + \left(\frac{z_t}{2}\right)(\tilde{h}_t + \hat{h}_t)$$

1.8 Step Down Network

$$y = \sigma(W^T h_t)$$

1.9 Error Calculation

$$E = \frac{1}{2}(\hat{y}_t - y_t)^2$$

2 Backpropagation of Weights at Time t

2.1 Update Gate Present Weights

$$\begin{aligned} \frac{\delta E_t}{\delta W_z} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\delta z_t^o} \cdot \frac{\delta z_t^o}{\delta z_t^i} \cdot \frac{\delta z_t^i}{\delta W_z} \\ \frac{\delta E_t}{\delta W_z} &= x_t [(W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{\tilde{h}_t}{2} + \frac{\hat{h}_t}{2} - h_{t-1}) \odot (\sigma(W_z x_t + U_z x_t) \\ &\quad \odot (1 - \sigma(W_z x_t + U_z h_{t-1})))^T] \end{aligned}$$

2.2 Update Gate Memory Weights

$$\begin{aligned} \frac{\delta E_t}{\delta U_z} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\delta z_t^o} \cdot \frac{\delta z_t^o}{\delta z_t^i} \cdot \frac{\delta z_t^i}{\delta U_z} \\ \frac{\delta E_t}{\delta U_z} &= h_{t-1} [(W[(y - \hat{y}) \odot (W^T h_t) \odot (1 - (W^T h_t))]) \odot (\frac{\tilde{h}_t}{2} + \frac{\hat{h}_t}{2} - h_{t-1}) \odot (\sigma(W_z x_t + U_z x_t) \\ &\quad \odot (1 - \sigma(W_z x_t + U_z h_{t-1})))^T] \end{aligned}$$

2.3 Candidate Gate Present Weights

$$\begin{aligned}\frac{\delta E_t}{W_{\tilde{h}_t}} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\tilde{h}_t^o} \cdot \frac{\delta \tilde{h}_t^o}{\delta \tilde{h}_t^i} \cdot \frac{\delta \tilde{h}_t^i}{W_{\tilde{h}_t}} \\ \frac{\delta E_t}{W_{\tilde{h}_t}} &= x_t [(W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{z_t}{2}) \odot (1 - \tanh^2(W_{\tilde{h}_t} x_t + (U_{\tilde{h}_t} h_{t-1}) \odot r_t))]^T\end{aligned}$$

2.4 Candidate Gate Memory Weights

$$\begin{aligned}\frac{\delta E_t}{U_{\tilde{h}}} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\tilde{h}_t^o} \cdot \frac{\delta \tilde{h}_t^o}{\delta \tilde{h}_t^i} \cdot \frac{\delta \tilde{h}_t^i}{U_{\tilde{h}}} \\ \frac{\delta E_t}{U_{\tilde{h}_t}} &= (h_{t-1} \odot r_t) [(W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{z_t}{2}) \\ &\quad \odot (1 - \tanh^2(W_{\tilde{h}_t} x_t + (U_{\tilde{h}_t} h_{t-1}) \odot r_t))]^T\end{aligned}$$

2.5 Reset Gate Present Weights

$$\begin{aligned}\frac{\delta E_t}{\delta W_r} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\delta \tilde{h}_t^o} \cdot \frac{\delta \tilde{h}_t^o}{\delta \tilde{h}_t^i} \cdot \frac{\delta \tilde{h}_t^i}{\delta r_t^o} \cdot \frac{\delta r_t^o}{\delta r_t^i} \cdot \frac{\delta r_t^i}{\delta W_r} \\ \frac{\delta E_t}{\delta W_r} &= x_t [(W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{z_t}{2}) \odot (1 - \tanh^2(W_{\tilde{h}_t} x_t + (U_{\tilde{h}_t} h_{t-1}) \odot r_t)) \\ &\quad \odot (U_{\tilde{h}_t} h_{t-1}) \odot (\sigma(W_r x_t + U_r h_{t-1}) \odot (1 - \sigma(W_r x_t + U_r h_{t-1})))^T\end{aligned}$$

2.6 Reset Gate Memory Weights

$$\begin{aligned}\frac{\delta E_t}{\delta U_r} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\delta \tilde{h}_t^o} \cdot \frac{\delta \tilde{h}_t^o}{\delta \tilde{h}_t^i} \cdot \frac{\delta \tilde{h}_t^i}{\delta r_t^o} \cdot \frac{\delta r_t^o}{\delta r_t^i} \cdot \frac{\delta r_t^i}{\delta U_r} \\ \frac{\delta E_t}{\delta U_r} &= h_{t-1} [(W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{z_t}{2}) \odot (1 - \tanh^2(W_{\tilde{h}_t} x_t + (U_{\tilde{h}_t} h_{t-1}) \odot r_t)) \\ &\quad \odot (U_{\tilde{h}_t} h_{t-1}) \odot (\sigma(W_r x_t + U_r h_{t-1}) \odot (1 - \sigma(W_r x_t + U_r h_{t-1})))^T\end{aligned}$$

2.7 Attention Gate Weight

$$\begin{aligned}\frac{\delta E_t}{\delta W_a} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\delta \hat{h}_t} \cdot \frac{\delta \hat{h}_t}{\delta W_a} \\ \frac{\delta E_t}{\delta W_a} &= \gamma_t [(W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{z_t}{2}) \odot (1 - \tanh^2(W_a \gamma_t))]^T\end{aligned}$$

2.8 Learning Function Weight

$$\begin{aligned}\frac{\delta E_t}{W_\alpha} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\delta \hat{h}_t} \cdot \frac{\delta \hat{h}_t}{\delta \gamma_t} \cdot \frac{\delta \gamma_t}{\delta \alpha_t} \cdot \frac{\delta \alpha_t}{\delta W_\alpha} \\ \frac{\delta E_t}{\delta W_\alpha} &= h_{t-1}^T x_t [((W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{z_t}{2}) \odot (1 - \tanh^2(W_a \gamma_t)) \\ &\quad \odot \text{softmax}'(\alpha_t))]^T\end{aligned}$$

2.9 Step-down Network Weight

$$\begin{aligned}\frac{\delta E_t}{\delta W} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta W} \\ \frac{\delta E_t}{\delta W} &= h_t [(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]^T\end{aligned}$$

3 Backpropagation of State at time t

3.1 Past State of Update Gate

$$\begin{aligned}\frac{\delta E_t}{\delta h_{t-1}} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\delta z_t^o} \cdot \frac{\delta z_t^o}{\delta z_t^i} \cdot \frac{\delta z_t^i}{\delta h_{t-1}} \\ \frac{\delta E_t}{\delta h_{t-1}} &= U_z [(W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{\tilde{h}_t}{2} + \frac{\hat{h}_t}{2} - h_{t-1}) \odot (\sigma(W_z x_t + U_z x_t) \\ &\quad \odot (1 - \sigma(W_z x_t + U_z h_{t-1})))]\end{aligned}$$

3.2 Past State of Candidate Gate

$$\begin{aligned}\frac{\delta E_t}{h_{t-1}} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\tilde{h}_t^o} \cdot \frac{\delta \tilde{h}_t^o}{\delta \tilde{h}_t^i} \cdot \frac{\delta \tilde{h}_t^i}{h_{t-1}} \\ \frac{\delta E_t}{h_{t-1}} &= (U_{\tilde{h}_t} \odot r_t) [(W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{z_t}{2}) \\ &\quad \odot (1 - \tanh^2(W_{\tilde{h}_t} x_t + (U_{\tilde{h}_t} h_{t-1}) \odot r_t))]\end{aligned}$$

3.3 Past State of Reset Gate

$$\begin{aligned}\frac{\delta E_t}{\delta h_{t-1}} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\delta \tilde{h}_t^o} \cdot \frac{\delta \tilde{h}_t^o}{\delta \tilde{h}_t^i} \cdot \frac{\delta \tilde{h}_t^i}{\delta r_t^o} \cdot \frac{\delta r_t^o}{\delta r_t^i} \cdot \frac{\delta r_t^i}{\delta h_{t-1}} \\ \frac{\delta E_t}{\delta h_{t-1}} &= U_r [(W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{z_t}{2}) \odot (1 - \tanh^2(W_{\tilde{h}_t} x_t + (U_{\tilde{h}_t} h_{t-1}) \odot r_t)) \\ &\quad \odot (U_{\tilde{h}_t} h_{t-1}) \odot (\sigma(W_r x_t + U_r h_{t-1}) \odot (1 - \sigma(W_r x_t + U_r h_{t-1})))]\end{aligned}$$

3.4 Past State of Attention Gate

$$\begin{aligned}\frac{\delta E_t}{\delta h_{t-1}} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_t} \cdot \frac{\delta h_t}{\delta \hat{h}_t} \cdot \frac{\delta \hat{h}_t}{\delta h_{t-1}} \\ \frac{\delta E_t}{\delta h_{t-1}} &= W_{\alpha_t} x_t [(W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t))]) \odot (\frac{z_t}{2}) \odot (1 - \tanh^2(W_a \gamma_t)) \odot softmax'(\alpha_t)]\end{aligned}$$

3.5 Past State of State Vector

$$\begin{aligned}\frac{\delta E_t}{\delta h_{t-1}} &= \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_{t-1}} \\ \frac{\delta E_t}{\delta h_{t-1}} &= (W[(y - \hat{y}) \odot \sigma(W^T h_t) \odot (1 - \sigma(W^T h_t)) \odot (1 - z_t)]\end{aligned}$$

3.6 Backpropagation State Value

$$\frac{\delta E_t}{\delta h_{t-1}} = \frac{\delta E_t}{\delta z_t^i} \cdot \frac{\delta z_t^i}{\delta h_{t-1}} + \frac{\delta E_t}{\delta \tilde{h}_t^i} \cdot \frac{\delta \tilde{h}_t^i}{h_{t-1}} + \frac{\delta E_t}{\delta r_t^i} \cdot \frac{\delta r_t^i}{\delta h_{t-1}} + \frac{\delta E_t}{\delta \hat{h}_t} \cdot \frac{\delta \hat{h}_t}{\delta h_{t-1}} + \frac{\delta E_t}{\delta y_t} \cdot \frac{\delta y_t}{\delta h_{t-1}}$$

4 Parameter Update

4.1 Update Update Gate

$$\nabla_{W_z} E = \sum_{t=0}^T \frac{\delta E_t}{\delta W_z} \quad \nabla_{U_z} E = \sum_{t=0}^T \frac{\delta E_t}{\delta U_z}$$

4.2 Update Reset Gate

$$\nabla_{W_r} E = \sum_{t=0}^T \frac{\delta E_t}{\delta W_r} \quad \nabla_{U_r} E = \sum_{t=0}^T \frac{\delta E_t}{\delta U_r}$$

4.3 Update Candidate Gate

$$\nabla_{W_c} E = \sum_{t=0}^T \frac{\delta E_t}{\delta W_c} \quad \nabla_{U_c} E = \sum_{t=0}^T \frac{\delta E_t}{\delta U_c}$$

4.4 Update Attention Weight

$$\nabla_{W_a} E = \sum_{t=0}^T \frac{\delta E_t}{\delta W_a}$$

4.5 Update Learning Function Weight

$$\nabla_{W_\alpha} E = \sum_{t=0}^T \frac{\delta E_t}{\delta W_\alpha}$$

4.6 Optimization Technique: NADAM

Algorithm 1 Nesterov-accelerated adaptive moment estimation

- 1: $g_t \leftarrow \nabla_{\theta_{t-1}} f(\theta_{t-1})$
 - 2: $\hat{g}_t \leftarrow \frac{g_t}{1 - \prod_{i=1}^t \mu_i}$
 - 3: $m_t \leftarrow \mu m_{t-1} + (1 - \mu) g_t$ $\triangleright \mu = .99$
 - 4: $\hat{m}_t \leftarrow \frac{m_t}{1 - \prod_{i=1}^{t+1} \mu_i}$
 - 5: $n_t \leftarrow v n_{t-1} + (1 - v) g_t^2$ $\triangleright v = .9$
 - 6: $\hat{n}_t \leftarrow \frac{n_t}{1 - v^t}$
 - 7: $\bar{m}_t \leftarrow (1 - \mu_t) \hat{g}_t + \mu_{t+1} \hat{m}_t$
 - 8: $\theta_t \leftarrow \theta_{t-1} - \eta \frac{\bar{m}_t}{\sqrt{\hat{n}_t} + \varepsilon}$ $\triangleright \varepsilon = 10e^{-8}$
-

5 Activation Functions

5.1 Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

5.2 Hyperbolic Tangent

$$\tanh(x) = \frac{e^x + e^{-x}}{e^x - e^{-x}}$$

5.3 Softmax

$$\text{softmax}(x)_j = \frac{e^{x_j}}{\sum_{k=1}^K e^{x_k}} \quad \text{for } j=1, \dots, K$$