

02/11/23

## Deep learning

o Recap

o GRU

o LSTM

o Recap: BPTT for RNN

$$\begin{aligned} - \theta^{(r)} &= \theta^{(r-1)} - \eta \nabla_{\theta^{(r-1)}} R(\theta) \quad (3) \\ - R(\theta) &= \frac{1}{T} \sum_{t=1}^T d(y^{(t)}, \hat{y}^{(t)}) \quad (4) \end{aligned}$$

Recap:

$$\begin{aligned} z_m^{(t)} &= \sigma(x_m^T \cdot x^{(t)} + \beta_m^T \cdot z^{(t-1)}) \quad (1) \\ \hat{y}^{(t)} &= g_h(\beta_k^T \cdot z^{(t)}) \quad (2) \end{aligned}$$

-  $\frac{\partial R(\theta)}{\partial \gamma_{mp}}$  expression ;  $\gamma_{mp}$  is the parameter connecting the  $p$ th hidden node at time  $t-1$  to the  $m$ th

hidden node at time  $t$ .

$$- \frac{\partial R(\theta)}{\partial \gamma_{mp}} = \frac{\partial R(\theta)}{\partial y} \cdot \frac{\partial y}{\partial z_m} \cdot \frac{\partial z_m^{(t)}}{\partial \gamma_{mp}} \quad (5)$$

-  $\frac{\partial z_m^{(t)}}{\partial \gamma_{mp}}$  is now the expensive computation (6)

- BPTT (full) is expensive computationally

- BPTT also leads to vanishing/exploding gradients

- Truncated BPTT to overcome these issues.

o Gated Recurrent Unit (GRU):

$$z_m^{(t)} = s_m^{(t)} \underbrace{z_m^{(t-1)}}_{\text{update gate}} + (1 - s_m^{(t)}) \underbrace{z_m^{(t)}}_{\text{candidate update}} \quad (7)$$

↳ update  
gate

$$\tilde{z}_m^{(t)} = \sigma(x_m^T \cdot x^{(t)} + r_m^{(t)} \cdot y_m^T \cdot z^{(t-1)}) \quad (8)$$

$$\hat{z}_m^{(t)} = \tanh(\underline{\alpha}_m^T \underline{x}^{(t)} + \underline{r}_m^{(t)} \underline{\beta}_m^T \underline{z}^{(t-1)}) \quad (9)$$

↳ reset gate

$$\hat{y}_k^{(t)} = g_k(\underline{\beta}_k^T \underline{z}^{(t)}). \quad (10)$$

### o GRU vs RNN

- $s_m^{(t)}$  - allows us to weight or give importance to the past state  $\underline{z}_m^{(t-1)}$

- this in turn allows for paying attention in the long range.

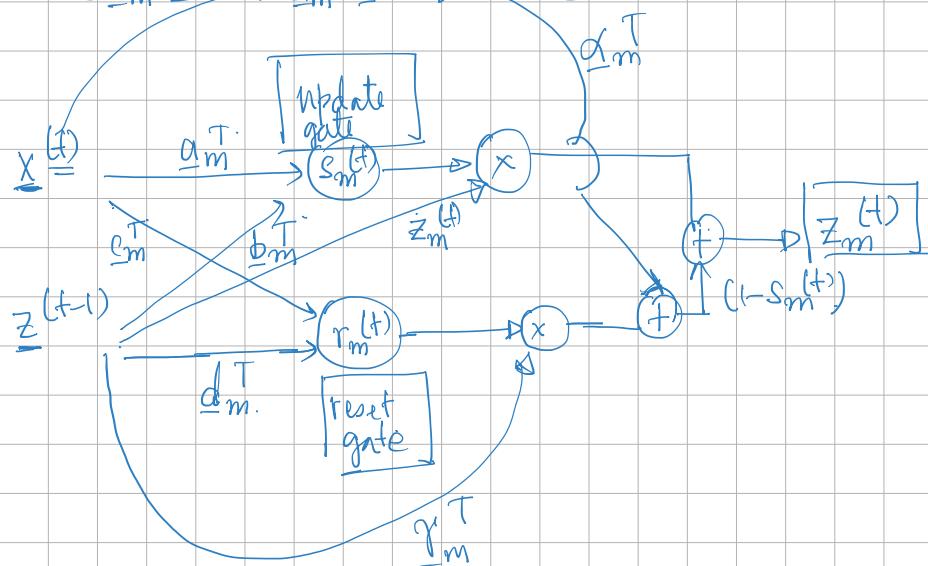
- $r_m^{(t)}$  - allows us to reset the past hidden state  $\underline{z}_m^{(t-1)}$  or pay lesser importance to it. This allows the model to pay attention to the short term or near term input.

$\therefore$  We now have both longterm and short term control knobs.

- The added parameters are the update gate  $s_m^{(t)}$  & the reset gate  $r_m^{(t)}$ .
- Both  $s_m^{(t)}$  and  $r_m^{(t)}$  take values in  $[0, 1]$ .

$$s_m^{(t)} = \sigma(\underline{\alpha}_m^T \underline{x}^{(t)} + \underline{b}_m^T \underline{z}^{(t-1)}) \quad (11)$$

$$r_m^{(t)} = \sigma(\underline{\beta}_m^T \underline{x}^{(t)} + \underline{d}_m^T \underline{z}^{(t-1)}) \quad (12)$$



$$\Theta = \{\underline{\alpha}, \underline{\beta}, \underline{g}, \underline{a}, \underline{b}, \underline{c}, \underline{d}\}$$

- Bias is a part of these parameter vectors  $\underline{a}_m, \underline{b}_m, \underline{c}_m, \underline{d}_m$

### o Long-Short Term Memory (LSTM)

RNN:

$$\underline{z}_m^{(t)} = \sigma(\underline{\alpha}_m^T \underline{x}^{(t)} + \underline{r}_m^{(t)} \underline{\beta}_m^T \underline{z}^{(t-1)})$$

