

## Exercise 3 - Recurrent Networks and NLP

Abdalla Arafa, 428044  
Ehsan Attar, 427214  
Sai Leela Poduru, 428272  
Prateek Rathod, 428396  
Shruti Shrivastava, 428364  
Erik Schwede, 428240

December 2023

### 3.1. Backpropagation through Time

- i)
- ii)
- iii)
- iv)

### 3.2. Gated recurrent units

- i)

Given:

$$u_t = \sigma(w.h_{t-1} + w.x_t) \quad (1)$$

$$s_t = w.(h_{t-1} + x_t) \quad (2)$$

$$h_t = u_t.h_{t-1} + (1 - u_t).s_t \quad (3)$$

First, let's differentiate equation (3) with respect to  $h_{t-1}$ :

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{u_t * h_{t-1} + (1 - u_t) * s_t}{\partial h_{t-1}} \quad (4)$$

On applying the chain rule :

$$\frac{\partial h_t}{\partial h_{t-1}} = u_t \cdot \frac{\partial h_{t-1}}{\partial h_{t-1}} + (1 - u_t) \cdot \frac{\partial s_t}{\partial h_{t-1}} \quad (5)$$

We have the value of s from equation 2:

$$(6)$$

$$s_t = w.(h_{t-1} + x_t) \quad (7)$$

$$\frac{\partial s_t}{\partial h_{t-1}} = w.\frac{\partial h_{t-1}}{\partial h_{t-1}} + 0 \quad (8)$$

$$= w \quad (9)$$

On substituting it back to our previous equation, we have:

$$\frac{\partial h_t}{\partial h_{t-1}} = u_t + (1 - u_t).w \quad (10)$$

On comparing it with the required form  $A_t . w + B_t$ , we get:

$$A_t = 1 - u_t \quad (11)$$

$$B_t = u_t \quad (12)$$

ii)

The long-term derivative can be written as:

$$\frac{\partial h_t}{\partial h_0} = \frac{\partial h_t}{\partial h_{t-1}} \cdot \frac{\partial h_{t-1}}{\partial h_{t-2}} \cdot \frac{\partial h_{t-2}}{\partial h_{t-3}} \cdots \frac{\partial h_1}{\partial h_0} \quad (13)$$

From part (i) we have found that, for every term:

$$\frac{\partial h_t}{\partial h_{t-1}} = (1 - u_t).w + u_t \quad (14)$$

$u_t$  is the output of a sigmoid function that has a value in the range of (0,1). Hence, the term remains non-zero. This helps to avoid the vanishing gradient problem, as no exponential decay term can cause the gradients to vanish over long sequences. When  $u_t$  is 1, it implies that the new hidden state retains the value of the previous hidden state and when  $u_t$  is 0, it ignores the information from the previous step.