

# LING530F: Deep Learning for Natural Language Processing (DL-NLP)

**Muhammad Abdul-Mageed**

muhammad.mageed@ubc.ca

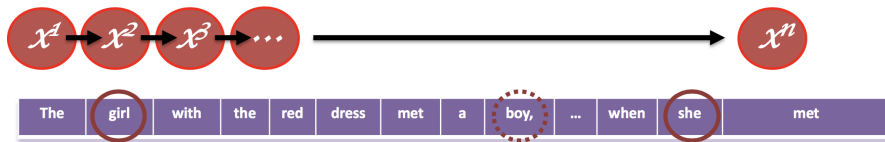
Natural Language Processing Lab

The University of British Columbia

# Table of Contents

- 1 Problems With Gradients
- 2 Gated Recurrent Units (GRU)
- 3 Long-Short Term Memory Networks (LSTMs)

# Vanishing and Exploding Gradients



**Figure:** A very long sequence, modelled with an RNN. Gradients can vanish and we will not know which gender a pronoun should be (male or female), or to which entity the pronoun refers (boy or girl)

## Gradient Problems

- Gradients can **explode**, in which case we can clip them.
- Gradients can also **vanish**, which is a more serious problem.

# Solving Long-Term Dependencies

## Solutions For Gradient Problems

- Long-Short Term Memory (**LSTM**) networks introduced to solve the problem of long-term dependencies
- Gated Recurrent Units (**GRU**) (Cho et al., 2014; Chung et al., 2014): Simplification of LSTMs.
- We will **introduce GRUs first**, as they are simpler
- **Notation modified from Andrew Ng**, rather than original papers, for pedagogical simplicity (**and elegance!**)

# Introducing a Memory Cell

augment network with a memory cell **C**

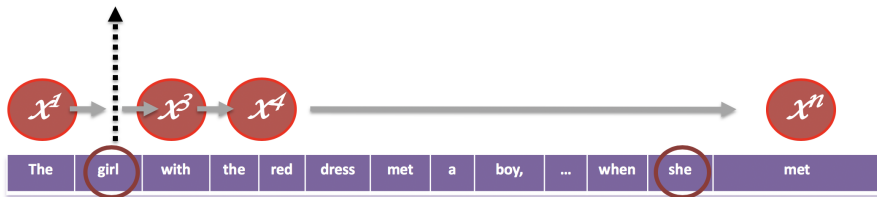


Figure: We will augment the network with a memory cell

## Gradient Problems

- The **memory cell** will help us retain information over long sequences
- For example, we can still know **we need pronoun she**, maintaining the **female gender** (and retaining the correct **reference** to "the girl")

## 1: Simple GRU

$$C_t = a_t$$

$$\tilde{C}_t = \tanh(W_c \cdot [C_{t-1}, x_t])$$

$$\Gamma_u = \sigma(W_u \cdot [C_{t-1}, x_t])$$

$$C_t = \Gamma_u * \tilde{C}_t + (1 - \Gamma_u) * C_{t-1}$$

## GRU Gates

- **Note:** Bias terms are dropped from equations
- $C_t$ : **Memory cell**.  $\tilde{C}_t$ : **New candidate memory cell for replacing  $C_t$** .
- $\Gamma_u$ : **Update gate**.

## 2: Simple GRU Update

$$C_t = \Gamma_u * \tilde{C}_t + (1 - \Gamma_u) * C_{t-1}$$

## GRU Gates

- $\Gamma_u$ : Result of a **a sigmoid** (between 0 and 1)
- When  $\Gamma_u$  is close to zero, we update very little  $\rightarrow$  we almost keep value of old memory cell  $C_{t-1}$

## 3: GRU With Relevance Gate

$$\tilde{C}_t = \tanh(W_c.[\Gamma_r * C_{t-1}, x_t])$$

$$\Gamma_u = \sigma(W_u.[C_{t-1}, x_t])$$

$$\Gamma_r = \sigma(W_r.[C_{t-1}, x_t])$$

$$C_t = \Gamma_u * \tilde{C}_t + (1 - \Gamma_u) * C_{t-1}$$

$$a_t = C_t$$

## GRU Gates

- $C_t$ : Memory cell.  $\tilde{C}_t$ : New candidate memory cell for replacing  $C_t$ .
- $\Gamma_u$ : Update gate.  $\Gamma_r$ : Relevance gate (aka reset gate)



## 4: GRU With Bias

$$\tilde{C}_t = \tanh(W_c \cdot [\Gamma_r * C_{t-1}, x_t] + b_c)$$

$$\Gamma_u = \sigma(W_u \cdot [C_{t-1}, x_t] + b_u)$$

$$\Gamma_r = \sigma(W_r \cdot [C_{t-1}, x_t] + b_r)$$

$$C_t = \Gamma_u * \tilde{C}_t + (1 - \Gamma_u) * C_{t-1}$$

## GRU Gates

- $C_t$ : Memory cell.  $\tilde{C}_t$ : New candidate memory cell for replacing  $C_t$ .
- $\Gamma_u$ : Update gate.  $\Gamma_r$ : Relevance gate (aka reset gate)

## 5: GRU With Bias

$$z_t = \sigma(W_z.[h_{t-1}, x_t])$$

$$r_t = \sigma(W_r.[h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W.[r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

## GRU Cells Notation Translation Table

- **z**: Update state  $\rightarrow \Gamma_u$
- **r**: Relevance state (*reset gate*)  $\rightarrow \Gamma_r$
- **$\tilde{h}$** : New candidate memory cell state  $\rightarrow \tilde{C}_t$
- **$h_t$** : GRU output  $\rightarrow C_t$

## 6: GRU With Bias

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z)$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t] + b_r)$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t] + b_h)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

## GRU Cells Notation Translation Table

- **z**: Update state  $\rightarrow \Gamma_u$
- **r**: Relevance state (*reset gate*)  $\rightarrow \Gamma_r$
- **$\tilde{h}$** : New candidate memory cell state  $\rightarrow \tilde{C}_t$
- **$h_t$** : GRU output  $\rightarrow C_t$

## 7: Recall: Full GRU

$$\tilde{C}_t = \tanh(W_c.[\Gamma_r * C_{t-1}, x_t])$$

$$\Gamma_u = \sigma(W_u.[C_{t-1}, x_t])$$

$$\Gamma_r = \sigma(W_r.[C_{t-1}, x_t])$$

$$C_t = \Gamma_u * \tilde{C}_t + (1 - \Gamma_u) * C_{t-1}$$

$$a_t = C_t$$

# Toward an LSTM

## Changes

- Parts in **red** will change!
- $a_t \neq C_t$  and so we will use  $a_t$
- We will also need to add new parts, namely a **forget gate** and **output gate** ...

## 8: Recall: Full GRU

$$\tilde{C}_t = \tanh(W_c.[\Gamma_r * C_{t-1}, x_t])$$

$$\Gamma_u = \sigma(W_u.[C_{t-1}, x_t])$$

$$\Gamma_r = \sigma(W_r.[C_{t-1}, x_t])$$

$$C_t = \Gamma_u * \tilde{C}_t + (1 - \Gamma_u) * C_{t-1}$$

$$a_t = C_t$$

# LSTM: New Candidate Cell $\tilde{C}_t$

## Updating $\tilde{C}_t$

- To acquire  $\tilde{C}_t$ , instead of  $C_{t-1}$ , we use the new  $a_{t-1}$  (since  $a_{t-1}$  is acquired differently than  $C_{t-1}$ , as we explain later)

## 9: New Candidate $\tilde{C}_t$

$$\tilde{C}_t = \tanh(W_c.[a_{t-1}, x_t])$$

# LSTM: Update and Forget Gates

## Changes

- We will **not** use an **relevance gate** ( $\Gamma_r$ )
- Instead of using one update gate  $\Gamma_u$ , we will use **two gates** to control cell content:  $\Gamma_u$  (**update gate**) and  $\Gamma_f$  (**forget gate**)
- **Forget gate** will give the new memory cell the option to keep or forget the **old cell** ( $C_{t-1}$ ), but just add to it via **update gate** ( $\Gamma_u$ )  
→ See  $C_t$  below

## 10: LSTM

$$\Gamma_u = \sigma(W_u \cdot [a_{t-1}, x_t])$$

$$\Gamma_f = \sigma(W_f \cdot [a_{t-1}, x_t])$$

$$C_t = \Gamma_u * \tilde{C}_t + \Gamma_f * C_{t-1}$$

## Output Gate

- As mentioned, we will use an **output gate** ( $\Gamma_o$ )
- **Output gate** will enable us to update our  $a_t$  via element-wise multiplication by  $\Gamma_o$

## 11: Output Gate

$$\Gamma_o = \sigma(W_o \cdot [a_{t-1}, x_t])$$

$$a_t = \Gamma_o * \tanh(C_t)$$



## 12: LSTM

$$\tilde{C}_t = \tanh(W_c \cdot [a_{t-1}, x_t])$$

$$\Gamma_u = \sigma(W_u \cdot [a_{t-1}, x_t])$$

$$\Gamma_f = \sigma(W_f \cdot [a_{t-1}, x_t])$$

$$\Gamma_o = \sigma(W_o \cdot [a_{t-1}, x_t])$$

$$C_t = \Gamma_u * \tilde{C}_t + \Gamma_f * C_{t-1}$$

$$a_t = \Gamma_o * \tanh(C_t)$$

## 13: LSTM

$$\tilde{C}_t = \tanh(W_c \cdot [a_{t-1}, x_t] + b_c)$$

$$\Gamma_u = \sigma(W_u \cdot [a_{t-1}, x_t] + b_u)$$

$$\Gamma_f = \sigma(W_f \cdot [a_{t-1}, x_t] + b_f)$$

$$\Gamma_o = \sigma(W_o \cdot [a_{t-1}, x_t] + b_o)$$

$$C_t = \Gamma_u * \tilde{C}_t + \Gamma_f * C_{t-1}$$

$$a_t = \Gamma_o * \tanh(C_t)$$

# LSTM Schematic Illustrated

$$\tilde{C}_t = \tanh(W_c \cdot [a_{t-1}, x_t])$$

$$\Gamma_u = \sigma(W_u \cdot [a_{t-1}, x_t])$$

$$\Gamma_f = \sigma(W_f \cdot [a_{t-1}, x_t])$$

$$\Gamma_o = \sigma(W_o \cdot [a_{t-1}, x_t])$$

$$C_t = \Gamma_u * \tilde{C}_t + \Gamma_f * C_{t-1}$$

$$a_t = \Gamma_o * \tanh(C_t)$$

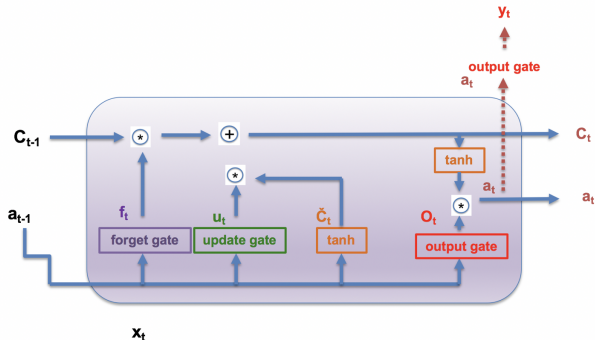
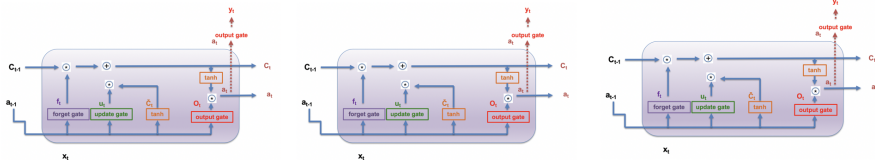


Figure: LSTM cell. [Inspired by Chris Olah and Andrew Ng]

# Stacking LSTM Cells



**Figure:** LSTM cell.s stacked. Note: Each cell will need new-indexing (not shown in the Figure)