

M2 Data Science - "Deep Learning 1" (IA-306)

Recurrent Neural Networks for Sequential Data

Geoffroy Peeters

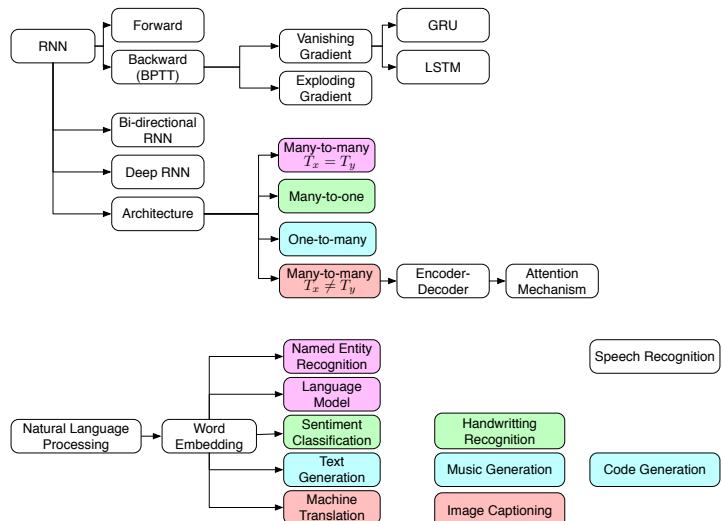
LTCI, Télécom Paris, IP Paris

2019 - 2020



Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 1

Overview



Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 3

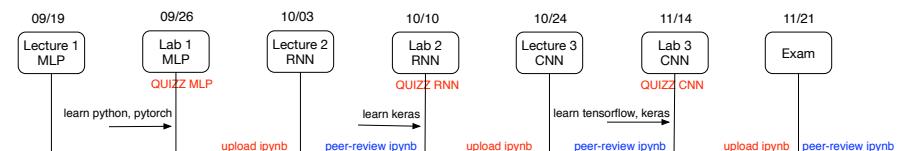
1. Three main types of Nets
 - 1.1 Multi Layers Perception (MLP) or Fully-Connected (FC)
 - 1.2 Recurrent Neural Networks (RNN)
 - 1.3 Convolutional Neural Networks (CNN)
2. Recurrent Neural Networks for Sequential Data
 - 2.1 What are sequential data ?
 - 2.2 Notations
 - 2.3 How can we process Sequential data with a Neural Network
 - 2.4 What is an RNN ?
 - 2.5 Forward Propagation
 3. Various types of sequential data
 - 3.1 Many-To-One, $T_x > 1, T_y = 1$
 - 3.2 Many-To-Many, $T_y = T_x$
 - 3.3 Many-To-Many, $T_y \neq T_x$
 - 3.4 One-To-Many, $T_x = 1, T_y > 1$
 4. Different architectures
 - 4.1 Bi-directional RNN
 - 4.2 Deep RNN
 - 4.3 Using 1D Convolution instead of RNN
 5. Back Propagation Through Time (BPTT)
 - 5.1 Forward pass
 - 5.2 Backward pass : 1) compute $\frac{\partial L}{\partial W_{ya}}$
 - 5.3 Backward pass : 2) compute $\frac{\partial L}{\partial W_{aa}}$
 - 5.4 Backward pass : 3) compute $\frac{\partial L}{\partial W_{ax}}$
 - 5.5 Vanishing and exploding gradients
 - 5.6 Dealing with the exploding gradient problem
 - 5.7 Dealing with the vanishing gradient problem
 6. Gated units (LSTM and GRU)
 - 6.1 Main ideas : Recurrent Neural Network (RNN)
 - 6.2 Main ideas : Long Short Term Memory Units (LSTM) LSTM regimes
 - 6.3 Main ideas : Gated Recurrent Unit (GRU)
 - 6.4 Recurrent Neural Network (RNN)
 - 6.5 Long Short Term Memory Units (LSTM)

- 6.6 Gated Recurrent Unit (GRU)
- 6.7 LSTM Example usage
7. Natural Language Processing
 - 7.1 Notations
 - 7.2 Language model
 - ① Training using a RNN
 - ② Testing using a RNN
 - ③ Sampling novel sequences using RNN
 - 7.3 One-hot-encoding
 - 7.4 Word embedding
 - How to get Word Embeddings
 - Use for Named Entity Recognition
 - Properties
 - Matrix
 - Learning using
 - Classic neural language model
 - Collobert and Weston model
 - Continuous Bag-Of-Words
 - Word2Vec/ Skip-gram model
 - Word2Vec/ Negative Sampling
 - Global vectors for word repres. (Glove)
 - 7.5 Application : Sentiment Classification using a simple model using RNN
 - 7.6 Application : Handwritten Character Recognition
 8. Sequence to sequence
 - 8.1 Introduction
 - 8.2 Machine Translation
 - ④ Decoding (finding the most likely sequence)
 - 8.3 Attention model
 - 8.4 Application : Image captioning
 - 8.5 Transformer
 9. Implementations in keras

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 2

Timetable/ Organization

Date	Type	Content	Teacher(s)	Location
September 19, 2019 PM	Lesson	MLP	G. Peeters	Télécom Paris (Amphi Estaunié)
September 26, 2019 PM	Lab	MLP	G. Peeters, A. Newson + others	Télécom Paris (C126, C127, C129, C130, C45)
October 3, 2019 PM	Lesson	RNN	G. Peeters	Télécom Paris (Amphi Estaunié)
October 10, 2019 PM	Lab	RNN	G. Peeters, A. Newson + others	Télécom Paris (C126, C127, C129, C130, C45)
October 24, 2019 PM	Lesson	CNN	A. Newson	Télécom Paris (Amphi Estaunié)
November 14, 2019	Lab	CNN	A. Newson, G. Peeters + others	Télécom Paris/Saclay (1A201, 1A207, 1A226, 1A252, 1A260)
November 21, 2019	Exam			Télécom Paris/Saclay (Amphi OB01)



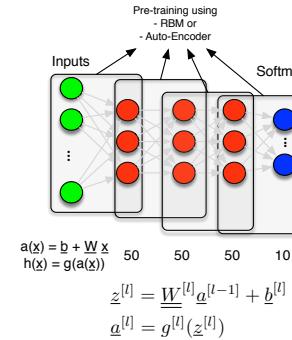
Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 4

Three main types of Nets

Three main types of Nets

Multi Layers Perceptron (MLP) or Fully-Connected (FC)

Multi Layers Perceptron (Fully Connected)

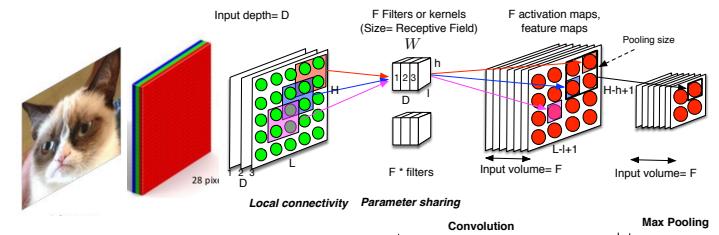


Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 6

Three main types of Nets

Convolutional Neural Networks (CNN)

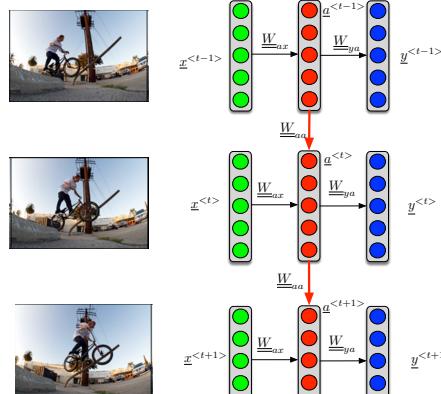
Convolutional Neural Network



Three main types of Nets

Recurrent Neural Networks (RNN)

Recurrent Neural Network



Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 7

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 8

Recurrent Neural Networks for Sequential Data



Recurrent Neural Networks for Sequential Data

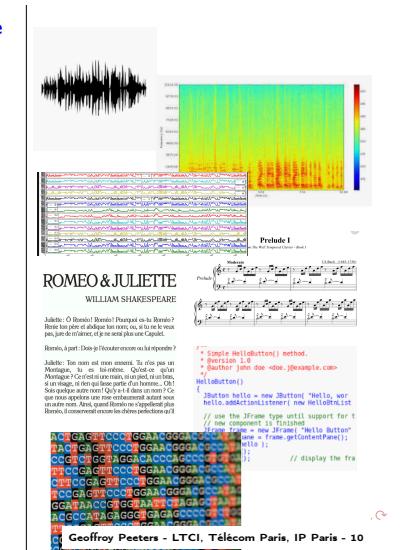
Notations

- Inputs :
 - $\underline{x}^{(t)}$: input vectors \underline{x} at time t (of dimension $n^{[0]}$)
 - $\{\underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(T_x)}\}$ the sequence of inputs of length T_x
- Outputs :
 - $\underline{y}^{(t)}$: input vectors \underline{y} at time t
 - $\{\underline{y}^{(1)}, \underline{y}^{(2)}, \dots, \underline{y}^{(T_y)}\}$ the sequence of outputs of length T_y
- T_x and T_y can be of different lengths
- $x^{(i)<t>}$ is the i^{th} example in the training-set
 - $T_x^{(i)}$ the length of the i^{th} input sequence
 - $T_y^{(i)}$ the length of the i^{th} output sequence

Recurrent Neural Networks for Sequential Data

What are sequential data ?

- **Sequential data are any type of data for which the order matters**
- Examples :
 - **Time series :**
 - Audio (sequence of sound pressures or sequence of Fourier transforms)
 - Sequence of musical notes
 - Video (sequence of images)
 - Sensors (heart-beat, plane altitude)
 - Stock market
 - **Ordered :** (but not with time)
 - Text/Words
 - Genes/ADN
- Common models :
 - Vector linear autoregression (VAR)
 - Markov model (Discrete-State HMMs)
 - Kalman filters (Continuous-State HMMs)
 - ...



Recurrent Neural Networks for Sequential Data

Notations

- Different types of output \underline{y} :
 - Continuous values $\underline{y} \in \mathbb{R}$ (Linear MSE Loss)
 - Binary-classes $\underline{y} \in \{0, 1\}$ (Sigmoid and Binary-Cross-Entropy)
 - Multi-classes $\underline{y} \in \{1 \dots K\}$ (Softmax and Cross-Entropy)
- Different types of input \underline{x}
 - Continuous values $\underline{x} \in \mathbb{R}$
 - Categorical values $\underline{x} \in \{1 \dots K\}$
- **Example** : categorical inputs/ binary outputs
 - "Named Entity Recognition"

$\{x\}$	$x^{<1>}$	$x^{<2>}$...	$x^{<7>}$		
Emmanuel Macron	is	the	president	of	France	
$\{y\}$	$y^{<1>}$	$y^{<2>}$...	$y^{<7>}$		
1	1	0	1	1	0	1

- **How do we process categorical values as input ?**
 - One-hot-encoding
 - Embedding

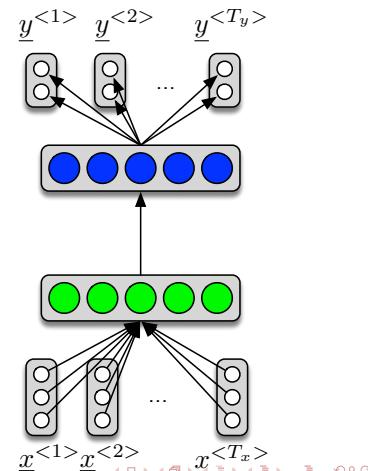
Recurrent Neural Networks for Sequential Data

How can we process Sequential data with a Neural Network

- Why Multi-Layer-Perceptron/Fully-connected networks are not appropriate?

- We could represent an input sequence $\{\underline{x}^{(1)} \dots \underline{x}^{(T_x)}\}$ as a large-vector of dimension $(n^{[0]} T_x)$
 - but (because T_x can be different for each input sequence) the dimension of this large-vector would change for each
 - zero-padding ?
 - still require a huge input dimension $(n^{[0]} T_x)$
 - the network will have to learn different weights for the same dimension $n^{[0]}$ at various time $\langle t \rangle$ in the sequence
 - no weight sharing !

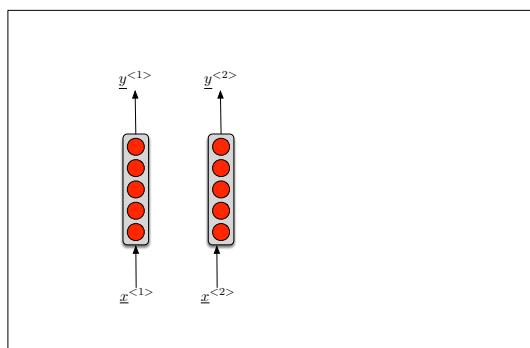
- **Solutions :**
 - Recurrent Neural Network
 - 1D Convolutional Neural Network
(convolution only over time)



Geoffrey Peeters - I-TCI, Télécom Paris, IP Paris - 1

Recurrent Neural Networks for Sequential Data

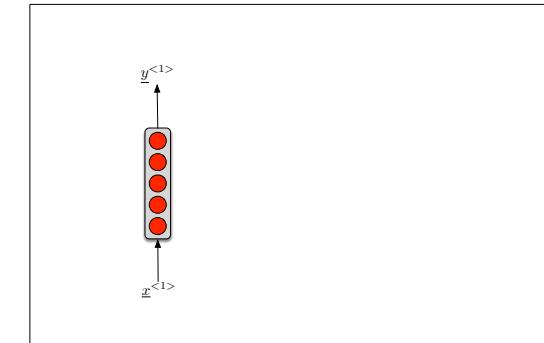
What is an RNN?



A set of small, light-blue navigation icons typically used in Beamer presentations for navigating between slides.

Recurrent Neural Networks for Sequential Data

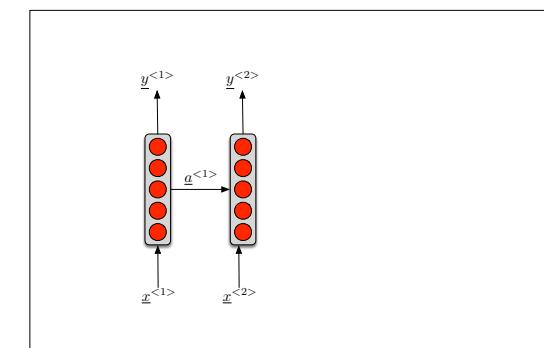
What is an RNN?



Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 14

Recurrent Neural Networks for Sequential Data

What is an RNN?

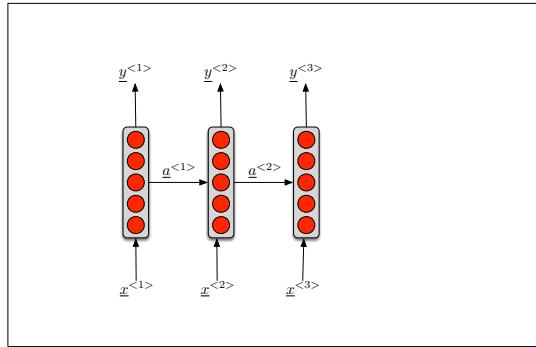


◀ □ ▶ ◀ 📂 ▶ ◀ 🔍 ▶ ◀ 🔍 ▶ 🔍 🔍 🔍

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 16

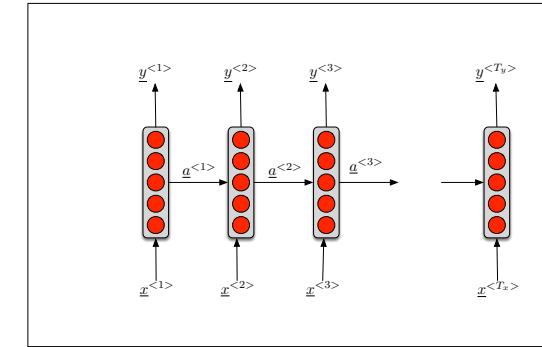
Recurrent Neural Networks for Sequential Data

What is an RNN ?



Recurrent Neural Networks for Sequential Data

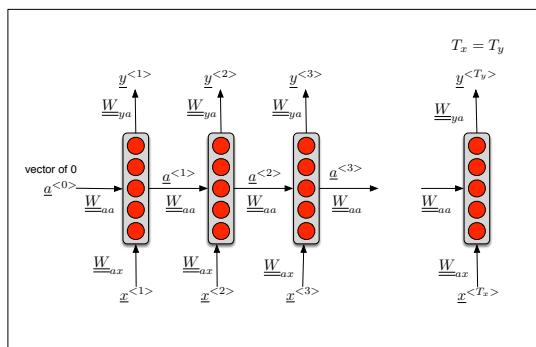
What is an RNN ?



- At each time step $< t >$, the RNN passes its activation $a^{<t>}$ to the next time step $< t + 1 >$

Recurrent Neural Networks for Sequential Data

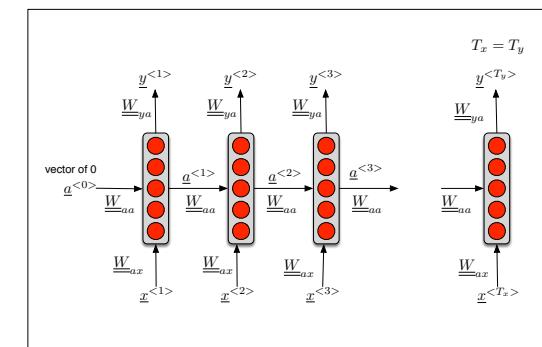
What is an RNN ?



- Initialize the first activation $a^{<0>}$ as zero

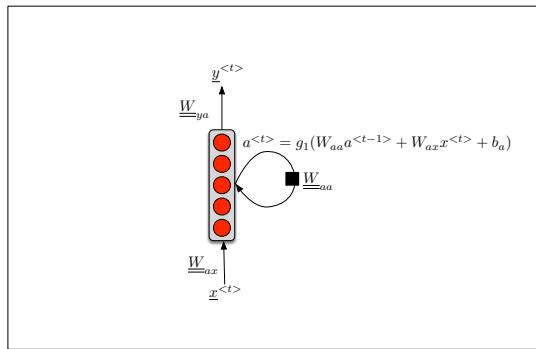
Recurrent Neural Networks for Sequential Data

What is an RNN ?



- \underline{W}_{ax} ($x^{<t>} \rightarrow a^{<t>}$) is the same for every time step
- \underline{W}_{aa} ($a^{<t-1>} \rightarrow a^{<t>}$) is the same for every time step
- \underline{W}_{ya} ($a^{<t>} \rightarrow y^{<t>}$) is the same for every time step

Other possible representation



- \underline{W}_{ax} ($x^{(t-1)} \rightarrow a^{(t)}$) is the same for every time step
 - \underline{W}_{aa} ($a^{(t-1)} \rightarrow a^{(t)}$) is the same for every time step
 - \underline{W}_{ya} ($a^{(t)} \rightarrow y^{(t)}$) is the same for every time step

Forward Propagation

- Forward propagation :

$$a^{<0\rangle} =$$

$$a^{<1>} = g_1 \left(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a \right)$$

$$\hat{y}^{(1)} = g_2 \left(W_{ya} a^{(1)} + b_y \right)$$

• • •

$$a^{(t)} = g_1 \left(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a \right)$$

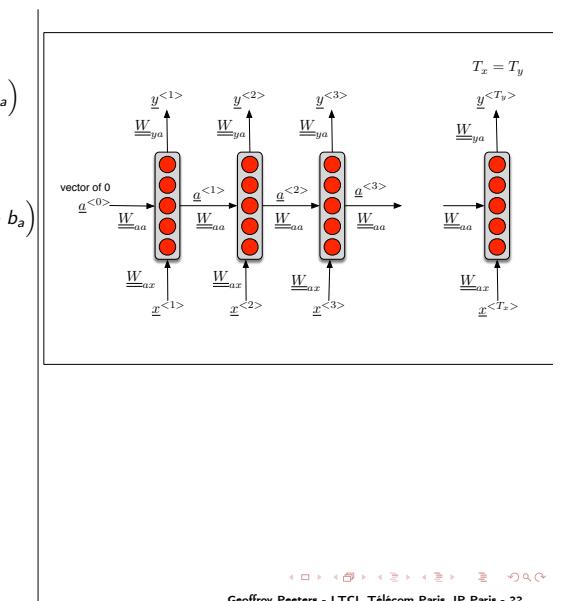
$$\hat{y}^{(t)} = g_2 \left(W_{ya} a^{(t)} + b_y \right)$$

- **Compact notations :**

$$\begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix} \left[\begin{array}{c} a^{(t-1)} \\ \vdots \\ \vdots \\ (n_a, m) \\ x^{(t)} \\ \vdots \\ (n_x, m) \end{array} \right]$$

$$a^{(t)} = g_1 \left(W_a [a^{(t-1)}; x^{(t)}] + b_a \right)$$

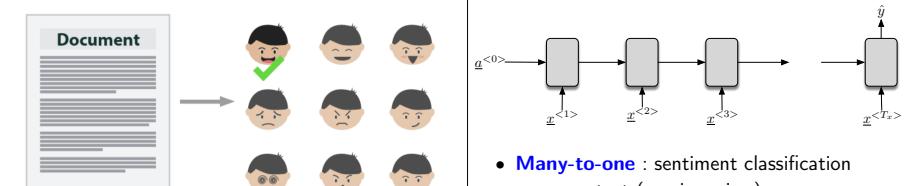
$$\hat{y}^{(t)} = g_2 \left(W_y a^{(t)} + b_y \right)$$



Various types of sequential data

Various types of sequential data/ Many-To-One $T_x > 1$, $T_y = 1$

<u>Input x</u>	<u>Output y</u>	Type	Examples
sequence $T_x > 1$	single $T_y = 1$	Many-To-One	Sentiment analysis, Video activity detection



- **Many-to-one** : sentiment classification

- x = text (movie review)
 - y = 0/1 or 1 · · · 5 rating



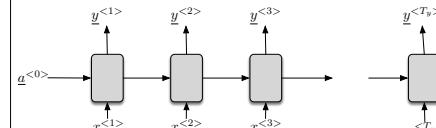
Various types of sequential data

Various types of sequential data/ Many-To-Many $T_y = T_x$

Input x	Output y	Type	Examples
sequence $T_x > 1$	sequence $T_y = T_x$	Many-To-Many	Named entity recognition

In 1917, **Einstein** applied the general theory of relativity to model the large-scale structure of the universe. He was visiting the **United States** when **Woodrow Wilson** came to power in 1913 and did not go back to **Germany**, where he had been a professor at the **Berlin Academy of Sciences**. He settled in the **U.S.**, becoming an American citizen in 1910. On the eve of World War II, he endorsed a letter to President **Franklin D. Roosevelt** alerting him to the potential development of "extremely powerful bombs of a new type" and recommending that the **US** begin similar research. This eventually led to what would become the **Manhattan Project**. Einstein supported defending the Allied forces, but largely denounced the use of the new discovery of nuclear fission as a weapon. Later, with the British philosopher **Sir Bertrand Russell**, Einstein signed the **Russell-Einstein Manifesto**, which highlighted the danger of nuclear weapons. Einstein was affiliated with the **Institute for Advanced Study** in Princeton, New Jersey, until his death in 1955.

Tag colours:
LOCATION TIME PERSON ORGANIZATION MONEY PERCENT DATE



- **Many-to-many** : $T_x = T_y$

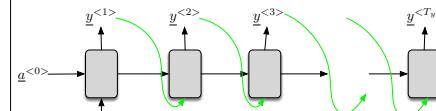
◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾
Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 25

Various types of sequential data

Various types of sequential data/ One-To-Many $T_x = 1, T_y > 1$

Input x	Output y	Type	Examples
single $T_x = 1$	sequence, $T_y > 1$	One-To-Many	Text generation, music generation

Nationalism and decision for the majority of Arab countries' capitulation was grounded by the Irish language by [[John Clancy]], [[An Imperial Japanese Revolt]], associated with Guangzhou's sovereignty. His generals were the powerful ruler of the Portuguese in the [[Protestant Immigrants]], which could be said to be directly in Cantonese Communication, which followed a ceremony and set inspired prison, training. The emperor travelled back to [[Antioch, Perth, October 25|21]] to note, the Kingdom of France, the United States, and the United Kingdom, the United Kingdom in western ([Scotland]), near Italy to the conquest of India with the conflict. Copyright was the succession of independence in the slip of Syrian influence that was a famous German movement based on a more popular servile, non-doctrinal and sexual power post. Many governments recognize the military housing of the [[Civil Liberalization and Infantry Resolution 265 National Party in Hungary]], that is synonymous to be the [[Civil Liberalization and Infantry Resolution]] (<http://www.hanabush@yahoo.com/guadianis>) (<http://7754800764175193e89.htm>) official economics Adjoint for the Nazism, Montgomery was swear to advance to the resources for those Socialism's rule, was starting to signing a major tripod of aid exile.]])



- **One-to-many** : text, code, music generation

AI-music, 192x96x576, 193 epochs, 805 seconds, 7 training sets



◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾
Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 27

Various types of sequential data

Various types of sequential data/ Many-To-Many $T_y \neq T_x$

Input x	Output y	Type	Examples
sequence $T_x > 1$	sequence $T_y > 1, T_y \neq T_x$	Many-To-Many	Automatic Speech Recognition, Machine translation



- **Many-to-many** : $T_x \neq T_y$: Machine translation

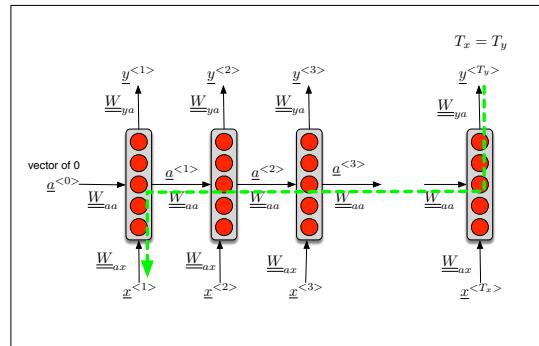
- Encoder/Decoder
- Attention mechanism
- Transformer architecture

◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾
Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 26

Different architectures

Different architectures

Bi-directional RNN



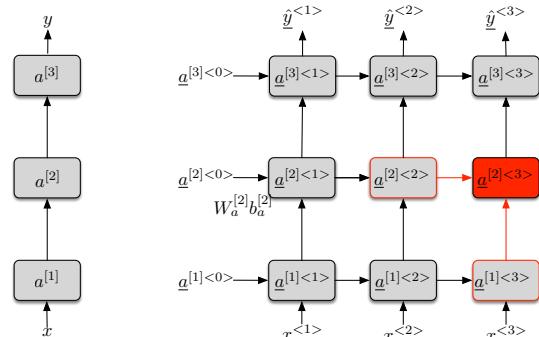
- RNN scans through the data from left to right \Rightarrow can do skip connections **but only in the past**
- Question : is "Apple" a named entity ?
 - "Apple stores sale the new **iPhone.**"
 - "Apple juice is a refreshing and healthy **beverages.**"
- \Rightarrow we cannot predict since we only depend on the past



Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 29

Different architectures

Deep RNN



- For the first layer, the inputs of the cell at time t are
 - the input x at the current time $\langle t \rangle$: $x^{(t)}$
 - the value of the cell at the previous time $\langle t-1 \rangle$: $a^{(t-1)}$
- For the layer $[l]$, the inputs of the cell at time $\langle t \rangle$ are
 - the value of the cell of the previous layer $[l-1]$ at the current time $\langle t \rangle$: $a^{[l-1]}_{\langle t \rangle}$
 - the value of the cell of the current layer $[l]$ at the previous time $\langle t-1 \rangle$: $a^{[l]}_{\langle t-1 \rangle}$

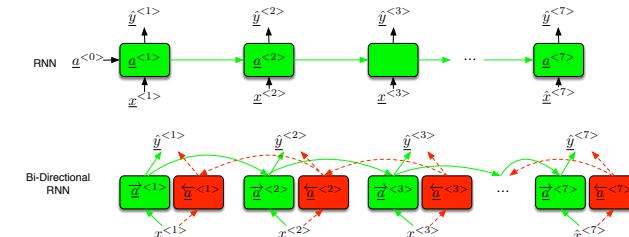
$$a^{[l]}_{\langle t \rangle} = g(W_a^{[l]}[a^{[l-1]}_{\langle t-1 \rangle}, a^{[l-1]}_{\langle t \rangle}] + b_a^{[l]})$$

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 31

Different architectures

Bi-directional RNN

- Solution ?
- Bi-directional RNN (BRNN)
- Question : is "Apple" a named entity ?
 - "Apple stores sale the new **iPhone.**"
 - "Apple juice is a refreshing and healthy **beverages.**"
- \Rightarrow we can predict since we now depend on the past and the future



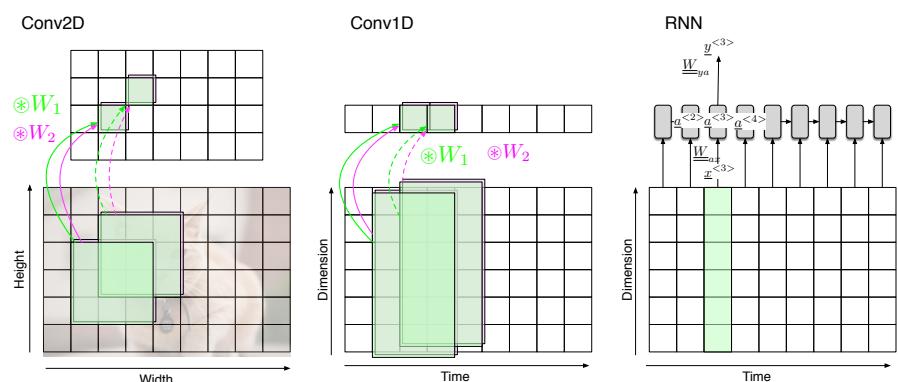
- The prediction $\hat{y}^{(t)}$ is done from the concatenation of
 - the **left-to-right** $\vec{a}^{(t)}$ and
 - the **right-to-left** $\overleftarrow{a}^{(t)}$ hidden states of RNN

$$\hat{y}^{(t)} = g(W_y[\vec{a}^{(t)}, \overleftarrow{a}^{(t)}] + b_y)$$

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 30

Different architectures

Using 1D Convolution instead of RNN



Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 32



Back Propagation Through Time (BPTT)

Backward pass : 1) compute $\frac{\partial \mathcal{L}}{\partial W_{ya}}$

- We measure how much varying W_{ya} affect $\mathcal{L}^{(t)}$.

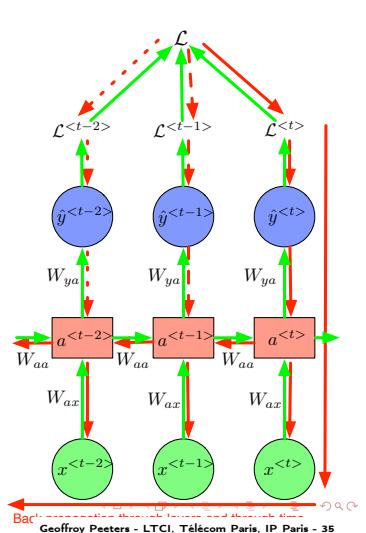
$$\frac{\partial \mathcal{L}}{\partial W_{ya}} = \sum_{t=1}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{ya}}$$

- We can use the chain rule.

$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_{ya}} = \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial W_{ya}}$$

- This is very simple since \hat{y}_t depends on W_{ya} on one place (indeed $a^{(t)}$ does not depend on W_{ya}).

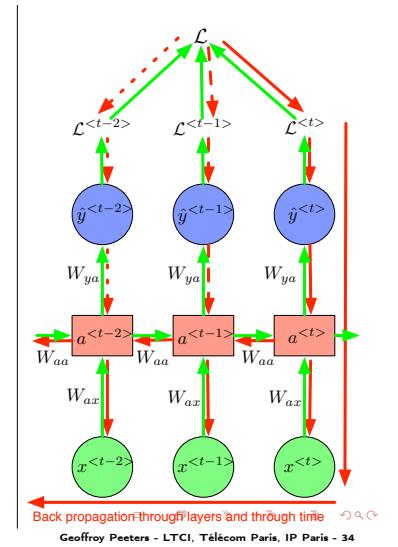
$$\hat{y}^{(t)} = g_y(W_{ya}a^{(t)} + b_y)$$



Back Propagation Through Time (BPTT)

Forward pass

- Compute $a^{(t)}, \hat{y}^{(t)}, \mathcal{L}^{(t)}, \mathcal{L}$
 $a^{(t)} = g_a(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a)$
 $\hat{y}_t = g_y(W_{ya}a^{(t)} + b_y)$
- Compute the Loss
 - Loss for time $t > 0$: $\mathcal{L}^{(t)}(y^{(t)}, \hat{y}^{(t)})$
 - $y^{(t)}$: true label
 - $\hat{y}^{(t)}$: predicted label
 - Total loss : $\mathcal{L} = \sum_t \mathcal{L}^{(t)}(y^{(t)}, \hat{y}^{(t)})$
- Backward pass**
 - compute $\frac{\partial \mathcal{L}}{\partial W_{ya}}, \frac{\partial \mathcal{L}}{\partial W_{ax}}, \frac{\partial \mathcal{L}}{\partial W_{aa}}, \frac{\partial \mathcal{L}}{\partial b_a}, \frac{\partial \mathcal{L}}{\partial b_y}$
 - All weights are shared across time steps !!!



Back Propagation Through Time (BPTT)

Backward pass : 2) compute $\frac{\partial \mathcal{L}}{\partial W_{aa}}$

- We measure how much varying W_{aa} affect $\mathcal{L}^{(t)}$.

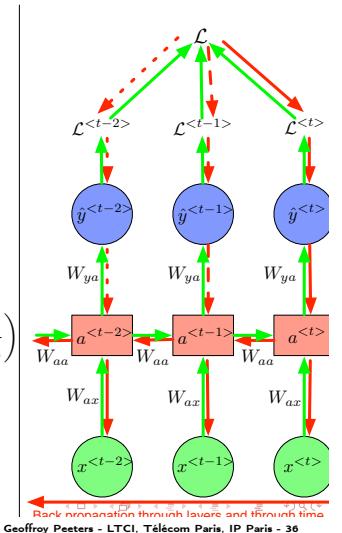
$$\frac{\partial \mathcal{L}}{\partial W_{aa}} = \sum_{t=1}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}}$$

- We cannot use the chain rule ($\frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}} = \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \frac{\partial a^{(t)}}{\partial W_{aa}}$) because a_t depends on W_{aa} also through $a^{(t-1)}$ which itself ...

$$a^{(t)} = g_a(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a)$$

- We use a formula of total derivative

$$\begin{aligned} \frac{d\mathcal{L}^{(t)}}{dW_{aa}} &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \left(\frac{\partial a^{(t)}}{\partial W_{aa}} + \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \frac{\partial a^{(t-1)}}{\partial W_{aa}} + \dots + \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \dots \frac{\partial a^{(0)}}{\partial W_{aa}} \right) \\ &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\hat{y}^{(t)}}{a^{(t)}} \sum_{k=0}^t \frac{\partial a^{(t)}}{\partial a^{(t-1)}} \dots \frac{\partial a^{(k+1)}}{\partial a^{(k)}} \frac{\partial a^{(k)}}{\partial W_{aa}} \\ &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial a^{(t)}} \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right) \frac{\partial a^{(k)}}{\partial W_{aa}} \end{aligned}$$



Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 34

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 35

Back Propagation Through Time (BPTT)

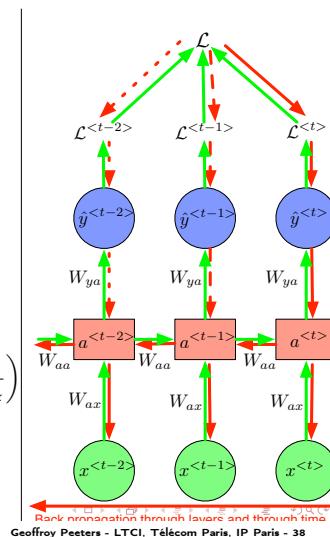
Backward pass : 3) compute $\frac{\partial \mathcal{L}}{\partial W_{ax}}$

- We measure how much varying W_{ax} affect $\mathcal{L}^{(t)}$.

$$\frac{\partial \mathcal{L}}{\partial W_{ax}} = \sum_{t=1}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{ax}}$$
 - It is also necessary to go backwards in time to calculate $\frac{\partial \mathcal{L}}{\partial W_{ax}}$
 - same situation as with W_{aa} :
 - $a^{(t)}$ depends on W_{ax} also through $a^{(t-1)}$ which itself ...

$$a^{(t)} = g_a(W_{ax}x^{(t)} + W_{aaa}a^{(t-1)} + b_a)$$
 - We use a formula of total derivative

$$\begin{aligned}
 & \frac{d\mathcal{L}^{(t)}}{dW_{ax}} \\
 &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial \mathbf{a}^{(t)}} \left(\frac{\partial \mathbf{a}^{(t)}}{\partial W_{ax}} + \frac{\partial \mathbf{a}^{(t)}}{\partial \mathbf{a}^{(t-1)}} \frac{\partial \mathbf{a}^{(t-1)}}{\partial W_{ax}} + \dots + \frac{\partial \mathbf{a}^{(t)}}{\partial \mathbf{a}^{(1)}} \frac{\partial \mathbf{a}^{(0)}}{\partial W_{ax}} \right) \\
 &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial \mathbf{a}^{(t)}} \sum_{k=0}^t \frac{\partial \mathbf{a}^{(t)}}{\partial \mathbf{a}^{(t-1)}} \dots \frac{\partial \mathbf{a}^{(k+1)}}{\partial \mathbf{a}^{(k)}} \frac{\partial \mathbf{a}^{(k)}}{\partial W_{ax}} \\
 &= \frac{\partial \mathcal{L}^{(t)}}{\partial \hat{y}^{(t)}} \frac{\partial \hat{y}^{(t)}}{\partial \mathbf{a}^{(t)}} \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial \mathbf{a}^{(i)}}{\partial \mathbf{a}^{(i-1)}} \right) \frac{\partial \mathbf{a}^{(k)}}{\partial W_{ax}}
 \end{aligned}$$

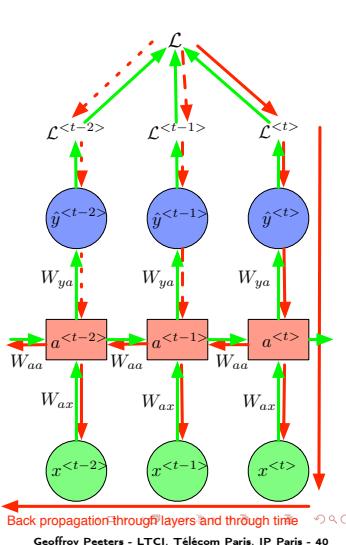


Back Propagation Through Time (BPTT)

Vanishing and exploding gradients

- To train an RNN we need to backpropagate through layers and through time
$$\frac{\partial \mathcal{L}}{\partial W_{aa}} = \sum_{t=0}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}}$$

$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}} \propto \sum_{k=0}^t \left(\prod_{i=k+1}^t \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right) \frac{\partial a^{(k)}}{\partial W_{aa}}$$
 - each term in the sum is the contribution of
 - a state at time k
 - to the gradient of the loss at time step t
 - the more steps between k and t , the more elements in this product
 - the values of these Jacobian matrices have particularly severe impact on the contributions of faraway steps



Back Propagation Through Time (BPTT)

Vanishing and exploding gradients

Vanishing gradient

- Long-term dependency
 - The **student**, which already followed 6 hours of lessons, **was** tired.
 - The **students**, which already followed 6 hours of lessons, **were** tired.
 - Very deep NN
 - very difficult to propagate back the gradient to affect the weights of the early Layers

Back Propagation Through Time (BPTT)

Vanishing and exploding gradients

$$\frac{\partial \mathcal{L}}{\partial W_{aa}} = \sum_{t=0}^T \frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}}$$

$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_{aa}} \propto \sum_{i=0}^t \left(\prod_{i'=\nu+1}^t \frac{\partial a^{(i')}}{\partial a^{(i-1)}} \right) \frac{\partial a^{(k)}}{\partial W_{aa}}$$

- Suppose only one hidden units, then a_i is a scalar and consequently $\frac{\partial a^{(i)}}{\partial a^{(i-1)}}$ is also a scalar
 - If $\left| \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right| < 1$ then the product goes to 0 exponentially fast
 - If $\left| \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right| > 1$ then the product goes to infinity exponentially fast
 - Vanishing gradients
 - $\left| \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right| < 1$
 - contributions from faraway steps vanish and don't affect the training
 - difficult to learn long-range dependencies
 - Exploding gradients
 - $\left| \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \right| > 1$
 - make the learning process unstable
 - gradient could even become NaN

Back Propagation Through Time (BPTT)

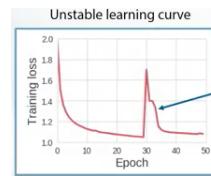
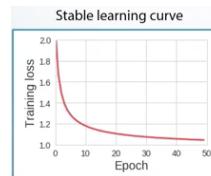
Dealing with the exploding gradient problem

- Solution 1 : gradient clipping**

- gradient $d\theta = \frac{\partial \mathcal{L}}{\partial \theta}$, where θ are all the parameters of the network
- if $\|d\theta\| > \text{threshold}$

$$d\theta \leftarrow \frac{\text{threshold}}{\|d\theta\|} d\theta$$

- clipping doesn't change the direction of the gradient but change its length
- we can clip only the norm of the part which causes the problem
- we choose the threshold manually : start with a large threshold and then reduce it



< □ > < ⏪ > < ⏴ > < ⏵ > < ⏹ > 🔍 🔍 🔍

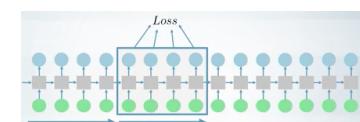
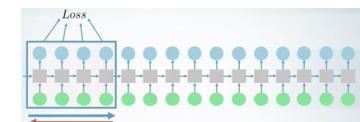
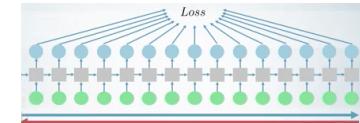
Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 45

Back Propagation Through Time (BPTT)

Dealing with the exploding gradient problem

- Solution 2 : truncated BPTT**

- Training very long sequences
 - time consuming !
 - exploding gradients !
- Truncated BPTT :
 - run forward and backward passes through the chunks of the sequence (instead of the whole sequence)
- Forward
 - We carry hidden states forward in time forever
- Backward
 - only backpropagate in the chunks (small number of steps)
- Much faster but dependencies longer than the chunk size don't affect the training (at least they still work at Forward pass)



< □ > < ⏪ > < ⏴ > < ⏵ > < ⏹ > 🔍 🔍 🔍

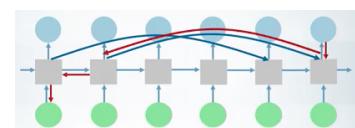
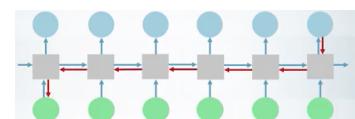
Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 46

Back Propagation Through Time (BPTT)

Dealing with the vanishing gradient problem

- Solution 2 : use skip-connections**

- in RNN : because at each step we multiply the Jacobian matrices → vanishing gradient
 - we cannot learn long-term dependencies
- Solution : add short-cuts between hidden states that are separated by more than one time step
 - are usual connections (with their own parameter matrices)
 - create much shorter paths between far away time-steps
- Backpropagation through the shortcuts : the gradients vanish slower
 - we can learn long-term dependencies
- not exclusive to RNN (see ResNet)



< □ > < ⏪ > < ⏴ > < ⏵ > < ⏹ > 🔍 🔍 🔍

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 49

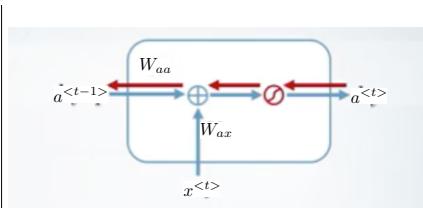
Gated units (LSTM and GRU)

Gated units (LSTM and GRU)

Main ideas : Recurrent Neural Network (RNN)

$$a^{(t)} = g(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a)$$

- non-linearities and/or multiplication creates the vanishing gradient problem
- Solution ?
 - create a short-way for backpropagation without any non-linearities or multiplication



Gated units (LSTM and GRU)

Main ideas : Long Short Term Memory Units (LSTM)

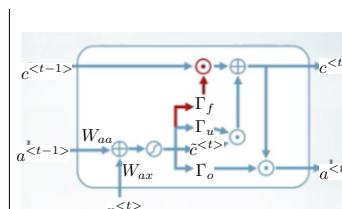
Forget gate LSTM :

Forget gate

$$\Gamma_f = \sigma(V_f x^{(t)} + W_f a^{(t-1)} + b_f)$$

$$c^{(t)} = \Gamma_f \odot c^{(t-1)} + \Gamma_u \odot \tilde{c}^{(t)}$$

$$a^{(t)} = \Gamma_o \odot f(c^{(t)})$$



- We now have a multiplication on the short-way

- $\frac{\partial c^{(t)}}{\partial c^{(t-1)}} = \text{diag}(\Gamma_f)$

- with Γ_f the forget gate : $\Gamma_f = \sigma(V_f x^{(t)} + W_f a^{(t-1)} + b_f)$

- Γ_f takes value between 0 and 1

- ⇒ set a high initial value for b_f

- at the beginning of the training, the LSTM doesn't forget and can find long-term dependencies in the data

- LSTM are very nice but has four times more parameters than RNN

Gated units (LSTM and GRU)

Main ideas : Long Short Term Memory Units (LSTM)

- Add a **memory cell**

- the same dimension as hidden layer

- Simplified LSTM** (no possibility to erase anything) :

Candidate cell value

$$\tilde{c}^{(t)} = g(V_g x^{(t)} + W_g a^{(t-1)} + b_g)$$

Update gate

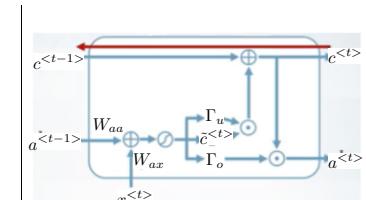
$$\Gamma_u = \sigma(V_u x^{(t)} + W_u a^{(t-1)} + b_u)$$

Output gate

$$\Gamma_o = \sigma(V_o x^{(t)} + W_o a^{(t-1)} + b_o)$$

$$c^{(t)} = c^{(t-1)} + \Gamma_u \odot \tilde{c}^{(t)}$$

$$a^{(t)} = \Gamma_o \odot f(c^{(t)})$$



- no non-linearity or multiplication in the memory cell-path

- $\frac{\partial c^{(t)}}{\partial c^{(t-1)}} = \text{diag}(1)$

- no vanishing gradients!!!

Gated units (LSTM and GRU)

Main ideas : Gated Recurrent Unit (GRU)

- Fewer gates**

Relevance gate

$$\Gamma_r = \sigma(V_r x^{(t)} + W_r a^{(t-1)} + b_r)$$

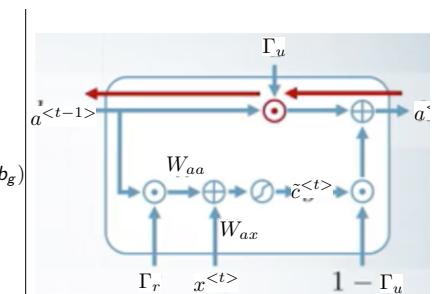
Candidate cell value

$$\tilde{c}^{(t)} = g(V_g x^{(t)} + W_g (a^{(t-1)} \cdot \Gamma_r) + b_g)$$

Update gate

$$\Gamma_u = \sigma(V_u x^{(t)} + W_u a^{(t-1)} + b_u)$$

$$a^{(t)} = (1 - \Gamma_u) \cdot \tilde{c}^{(t)} + \Gamma_u \cdot a^{(t-1)}$$



- Vanishing gradient problem ?**

- similar as for LSTM

$$\frac{\partial a^{(t)}}{\partial a^{(t-1)}} = \text{diag}(1 - \Gamma_u) \cdot \frac{\partial \tilde{c}^{(t)}}{\partial a^{(t-1)}} + \text{diag}(\Gamma_u)$$

- ⇒ high initial value for b_u

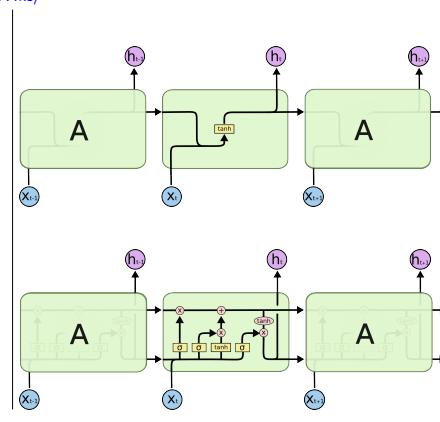
Gated units (LSTM and GRU)

LSTM Example usage

Based on <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Instead of having a single neural network layer, there are four, interacting in a very special way.

- Lines carry a vector
- A recurrent net transmits its hidden states
- LSTM introduces a second channel : the cell state
- It acts as a memory
- Gates control the memory

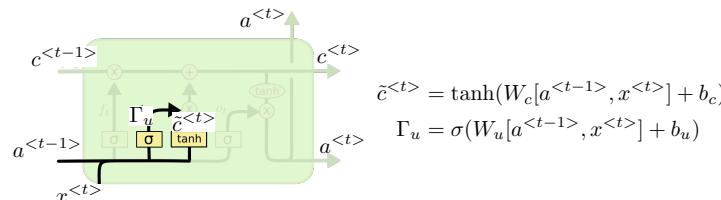


Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 61

Gated units (LSTM and GRU)

LSTM Example usage

- What should be taken into account ?



Action

- Create the update $\tilde{c}^{<t>}$ of the cell state
- and its contribution Γ_u (the update gate with a sigmoid activation)

Intuition for language modeling

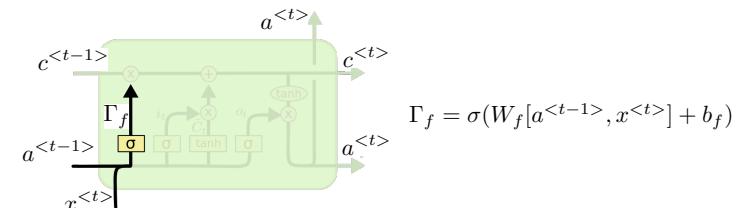
- Add the gender of the new subject to the cell state,
- to replace the old one we're forgetting. appears

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 63

Gated units (LSTM and GRU)

LSTM Example usage

- What should be forgotten from the previous cell state ?



Action

- The sigmoid (forget gate) answers for each component :
 - 1 : to keep it,
 - 0 to forget it, or
 - a value in-between to mitigate its influence

Intuition for language modeling

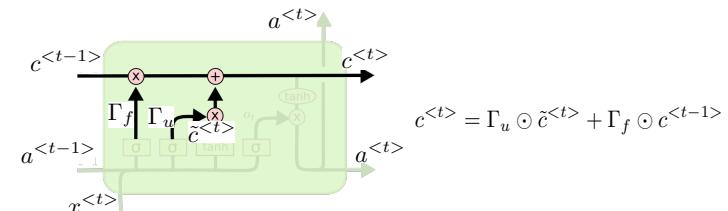
- The cell state might embed the gender of the present subject :
 - keep it to predict the correct pronouns,
 - or forget it, when a new subject appears

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 62

Gated units (LSTM and GRU)

LSTM Example usage

- Write the new state



Action

- Merge the old cell state modified by the forget gate
- with the new input

Intuition for language modeling

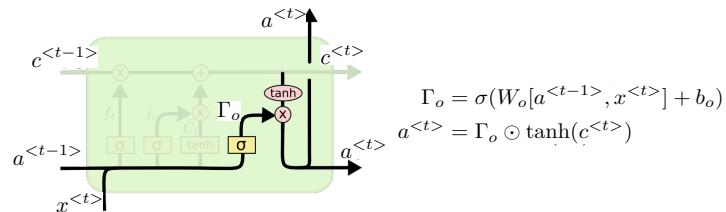
- Decide to drop the information about the old subject
- Refresh the memory appears

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 64

Gated units (LSTM and GRU)

LSTM Example usage

- Write the new hidden state



Action

- Decide what parts of the (filtered) cell state to output $a^{(t)}$
 - Compute the hidden state

Intuition for language modeling

- Since we just saw a subject,
 - output the relevant information for the future (gender, numnber) appears

Natural Language Processing

Language model

- **Automatic Speech Recognition :**
 - "This is an **example** of an homophone." or "This is an **egg-sample** of an homophone."
 - $p(\text{"example"}) = 5.7 \cdot 10^{-10}$
 - $p(\text{"egg - sample"}) = 3.2 \cdot 10^{-13}$
 - **Language model : $p(\text{sentence})$**
 - $p(y^{(1)}, y^{(2)}, \dots, y^{(T_y)})$
 - **How to build a language model ?**
 - need a large corpus of English text
 - tokenize the text
 - convert to one-hot-vector
 - + <EOS> (end of sentence)
 - + <UNK> (unknown word, very uncommon vocabulary) : Gif-sur-Yvette

Natural Language Processing



Natural Language Processing

Language model

$$\mathbf{w} \equiv w_1^L; w_1, w_2, \dots, w_L$$

- **Applications**
 - Automatic Speech Recognition, Machine Translation, OCR, ...

• Goal

- Estimate the (non-zero) probability of a word sequence for a given vocabulary

$$P(w_1^L) = P(w_1, w_2, \dots, w_L) = \prod_{i=1}^L P(w_i | w_1^{i-1}) \quad \forall i, w_i \in V$$

- with the **N-gram assumption**

$$P(w_1^L) = \prod_{i=1}^L P(w_i | w_{i-n+1}^{i-1}), \forall i, w_i \in V$$

- in the **recurrent way**

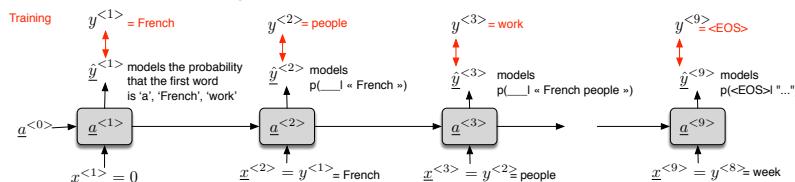
$$P(w_i | w_1^{i-1})$$

Natural Language Processing

Language model/ ① Training using a RNN

Algorithm :

- start with empty context, empty input word
- given the previous word as input $x^{(t)}$ and the context (in the hidden cell $a^{(t)}$) predict (maximize the probability to observe) the next word $y^{(t)}$ as output $\hat{y}^{(t)}$
- "French people work an average of 35 hours a week."



Training :

- output of the network (softmax) is $\hat{y}^{(t)}$ (probability of each word)
- ground-truth is $y^{(t)}$ (one-hot-encoding)
- minimize the loss

$$\mathcal{L}_t(\hat{y}^{(t)}, y^{(t)}) = - \sum_{c=1}^K y_c^{(t)} \log(\hat{y}_c^{(t)})$$

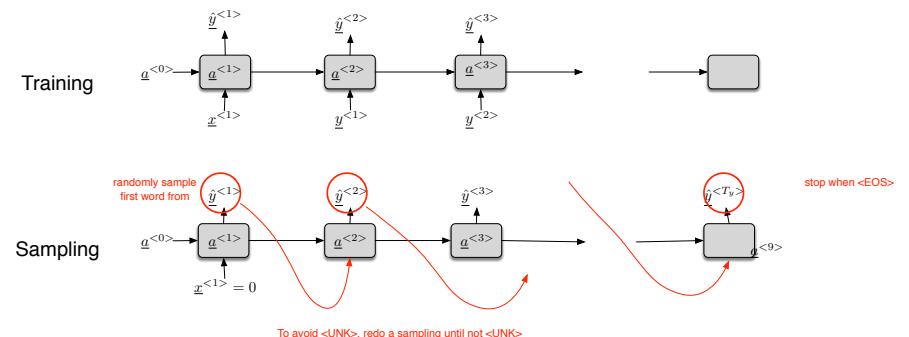
$$\mathcal{L} = \sum_t \mathcal{L}_t(\hat{y}^{(t)}, y^{(t)})$$

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 69

One-hot-encoding

Natural Language Processing

Language model/ ③ Sampling novel sequences using RNN



Word-level RNN

- Character-level RNN :
 - Vocabulary= ['a', 'b', 'c', ..., '.', ',', ';', 'A', 'B', 'C', ...]
- Examples :
 - Shakespeare, wikipedia, source-code generation
 - <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 71

Natural Language Processing

One-hot-encoding

- How to represent the individual words in a sentence ?
- What will be $x^{(t)}$?

- Construct the list of the N most used categories (words) in the corpus
 - = "vocabulary/dictionary"
 - N is usually in the range 10.000 words - 50.000.
 - add a <UNK> (unknown word) for the words that are not in vocabulary/dictionary
- Each word is now associated with an ID

1	Angela
...	Boris
...	chancelor
367	Emmanuel
...	France
...	Germany
4075	is
...	Johnson
6830	Macron
...	Merkel
...	of
...	president
...	prime-minister
...	the
10.000	UK

One-hot-encoding

- **One-hot-encoding :**
 - the one-hot vector of an ID is a vector filled with 0s, except for a 1 at the position associated with the ID
 - each word w is associated with a one-hot-vector vector o_{ID}
 - If $N = 10.000$, $x^{(t)}$ has dimension 10.000
- a one-hot encoding makes no assumption about word similarity
 - all words are equally different from each other

$\{x\}$	$x^{<1>}$	$x^{<2>}$	$x^{<3>}$	$x^{<4>}$	$x^{<5>}$	$x^{<6>}$	$x^{<7>}$
Emmanuel	0	0	0				
Macron	0	0	0				
is				1			
the					1		
president						1	
of							1
France							
1	0	0	0				
2	0	0	0				
...				
367	1	0	0				
...				
4075	0	0	1				
...				
6830	0	1	0				
...				
10.000	0	0	0				

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 74

One-hot-encoding

- **Weaknesses of one-hot-encoding**
 - it is very high-dimensional
 - vulnerability to overfitting, computationally expensive
 - all words are equally different from each other
 - the inner product between any two one-hot vectors is zero
 - as a consequence, we cannot generalize
- Example :
 - we have a trained a language model to complete the sentence
 - The president is on a **boat** ? " \Rightarrow "trip"
 - since there is no specific relationship between "boat" and "plane", the algorithm cannot complete the sentence
 - The president is on a **plane** ? " \Rightarrow ?



Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 75

Word embedding

Word embedding

- **Word Embedding** : learn a continuous representation of words
 - each word w is associated with a real-valued vector e_{ID}
 - if embedding size is 300, e_{5391} is a 300 dimensional vector
 - we would like the distance $\|e_{ID1} - e_{ID2}\|$ to reflect the meaningful similarities between words

	Man 5391	Woman 9853	King 4914	Queen 7157	Uncle 456	Aunt 6257	President 7124	Chancellor 3212	France 6789	Germany 1234	Boat 5923	Plane 8871
Gender	-1	1	-1	1	-1	1	-0,7	-0,3	0	0	0	0
Royal	0	0	1	1	0	0	0,5	0,5	0	0	0	0
Age	0	0	0,7	0,7	0,5	0,5	0,7	0,7	0	0	0	0
Country	0	0	0,5	0,5	0	0	0,5	0,2	1	1	0	0
Transportation	0	0	0	0	0	0	0	0	0	0	1	1
...

- In this representation "boat" and "plane" are quite similar
 - some of the features may differ but most features are similar
- Therefore the algorithm can find (by generalization) that
 - The president is on a **boat** ? " \Rightarrow "trip"
 - The president is on a **plane** ? " \Rightarrow "trip"



Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 77

Natural Language Processing

Word embedding/ How to get Word Embeddings

- ① Train from scratch a Word Embedding matrix
 - require a lot of data (1 to 100 Billion words)
 - ② Download from the internet
 - <https://fasttext.cc/docs/en/crawl-vectors.html>
 - has been already trained on 1 to 100 Billion words
 - perform **transfer learning** :
 - use directly for a task (possibly with a much smaller training set)
 - use after some fine-tuning the word embeddings with the new training set



Natural Language Processing

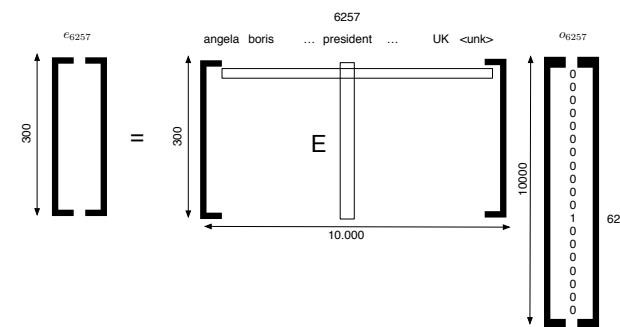
Word embedding/ Learning using

- Various existing methods to learn word embedding :
 - Classic neural language model (Bengio)
 - Collobert and Weston model
 - Word2Vec (Mikolov)
 - Continuous bag-of-words (CBOW)
 - Skip-gram
 - Glove model (Pennington)



Natural Language Processing

Word embedding/ Matrix



- \underline{o}_{6257} = one-hot vector (zero everywhere except 1 at position 6257)
 - dimension 10.000
 - $\underline{e}_{6257} = \underbrace{\underline{E}}_{(300)} \underbrace{\underline{o}_{6257}}_{(300,10.000)(10.000,1)}$ = embedding of \underline{o}_{6257}
 - $e_j = \underline{E} \cdot \underline{o}_j$ = embedding for word j
 - In practice, use look-up table (take the column vector of E corresponding to 6257)

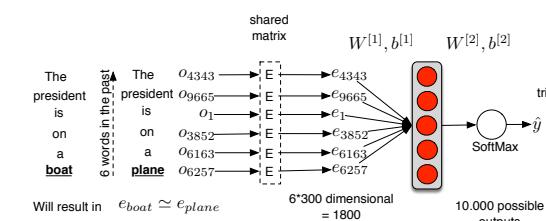
Gaffney Brothers - LTCI Telecom Basis IP Basis - 83

Natural Language Processing

Word embedding/ Classic neural language model

[Bengio et al., 2003 "A Neural probabilistic model"]

- **Method** : build a language model (learn to predict the following word) using a MLP
 - Use backpropagation to learn the parameters : $E, W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
 E is the **embedding matrix**



- **Example :**
 - "The president is on a boat trip_{target} to the United States."
 - **Other possible definition of the context**
 - last 4 words
 - 4 words on the left, 4 words on the right : "is on a boat ? to the United States"
 - last 1 words : "boat ?"
 - nearby 1 word (skip-gram) : "president ?"

◀ □ ▶ ⏪ ⏩ ⏴ ⏵ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 84

Natural Language Processing

Word embedding/ Collobert and Weston model

[Collobert and Weston, 2008 "A unified architecture for natural language processing: Deep neural networks with multitask learning."]

- Goal :
 - avoid having to compute the expensive softmax

- How ?

- different objective function
 - instead of the cross-entropy criterion of Bengio (maximizes the probability of the next word given the previous words),
 - train a network to output a higher score f_θ for a correct word sequence than for an incorrect one.
 - pairwise ranking criterion :

$$J_\theta = \sum_{x \in X} \sum_{w \in V} \max \{0, 1 - f_\theta(x) + f_\theta(x^{(w)})\}$$

- From the set of all possible windows X in the corpus,
 - sample **correct windows** x containing n words
- For each window x :
 - **produce a corrupted/incorrect version $x^{(w)}$** (by replacing the center word of x with another word w from V)

Application : Sentiment Classification

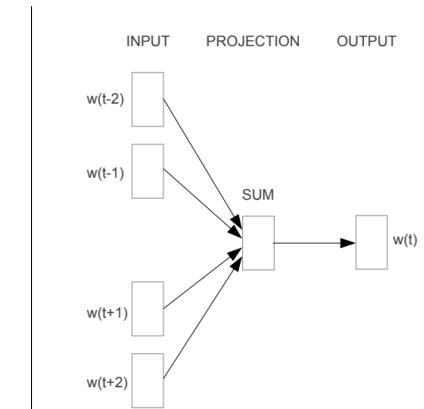
Natural Language Processing

Word embedding/ Continuous Bag-Of-Words

[Mikolov et al. 2013. "Efficient estimation of word representations in vector space."]

- **Mikolov/ CBOW** :

- predicts the current word based on the context
 - for target word $w(t)$, use both the n words before and after
 - named "continuous bag-of-words (CBOW)" ?
 - uses continuous representations, order is of no importance (SUM)



Natural Language Processing

Application : Sentiment Classification

x

y

This movie is fantastic ! I really like it because it is so good !



Not to my taste, will skip and watch another movie.



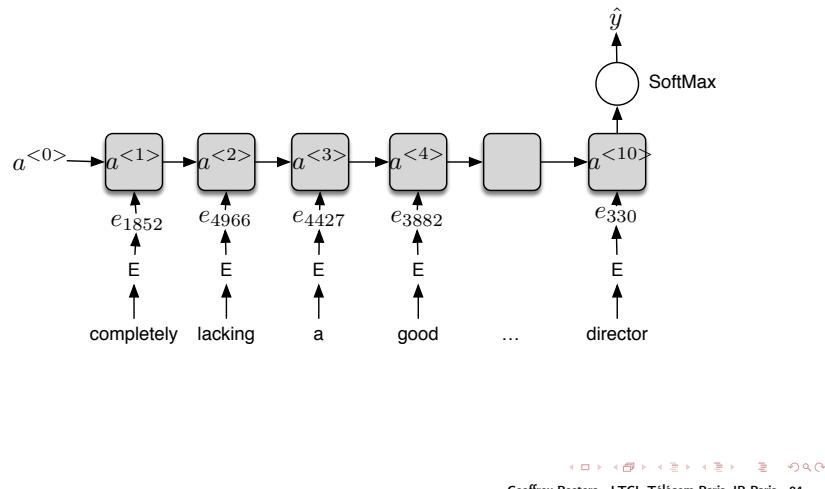
This movie was completely lacking a good script, good actors and a good director.



Natural Language Processing

Application : Sentiment Classification/ using RNN

- Many-to-one architecture

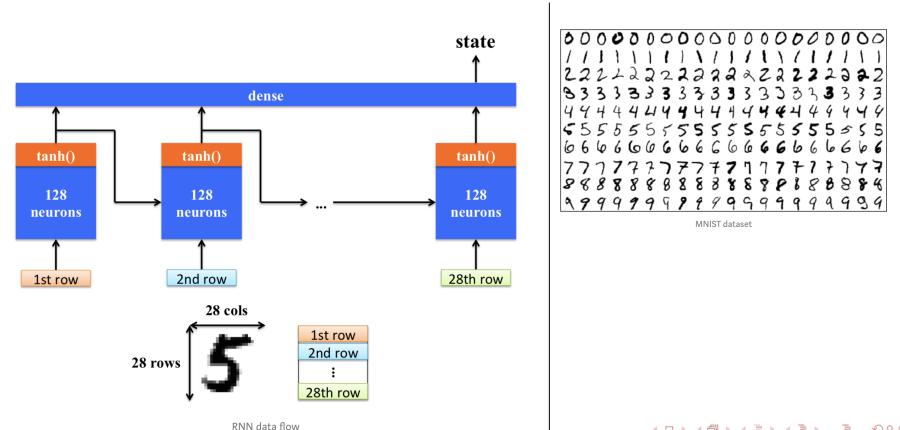


Natural Language Processing

Application : Handwritten Character Recognition

[Graves et al. 2008 "Unconstrained online handwriting recognition with recurrent neural networks"]
 [Graves et al. 2009 "Offline Handwriting Recognition with Multidimensional Recurrent Neural Networks"]

- source <https://medium.com/machine-learning-algorithms/mnist-using-recurrent-neural-network-2d070a5915a2>



0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6
 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9

MNIST dataset

Sequence to sequence

Sequence to sequence

Introduction

[Sutskever et al. 2014, "Sequence to sequence learning with neural networks"]

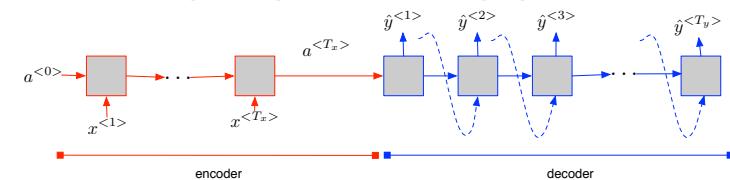
[Cho et al. 2014 "Learning phrase representations using rnn encoder-decoder for statistical machine translation"]

- What happens if $T_x \neq T_y$ (many-to-many architecture but with different lengths)
 - Example : Machine Translation

$\langle x \rangle$	$x^{<1>}$	$x^{<2>}$...		$x^{<6>}$	
	Macron	voyage	aux	Etats-Unis à	Noël	
$\langle y \rangle$	$y^{<1>}$	$y^{<2>}$...			$y^{<8>}$
	Macron	is	travelling	to	the	United-Statfor

- Solution :

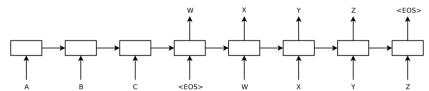
- Sequence to sequence [Sutskever] or Encoder-Decoder [Cho] architecture



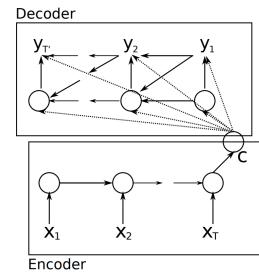
Sequence to sequence

Introduction

[Sutskever et al. 2014, "Sequence to sequence learning with neural networks"]



[Cho et al. 2014 "Learning phrase representations using rnn encoder-decoder for statistical machine translation"]

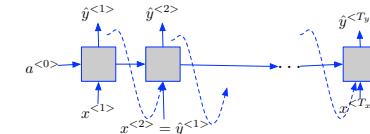


Sequence to sequence

Machine Translation

- **Language model :**

$$p(y^{(1)}, \dots, y^{(T_y)})$$

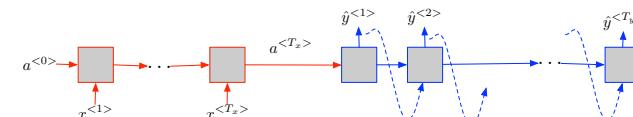


- **Machine Translation :**

$$\text{Conditional language model : } p(y^{(1)}, \dots, y^{(T_y)} | x^{(1)}, \dots, x^{(T_x)})$$

- French sentence : $x^{(1)}, \dots, x^{(T_x)}$

- English sentence : $y^{(1)}, \dots, y^{(T_y)}$



Attention model

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 98

Sequence to sequence

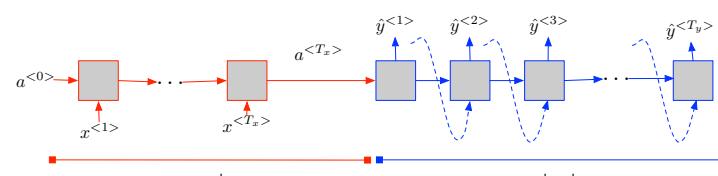
Attention model

[Bahdanau et al. 2014 "Neural machine translation by jointly learning to align and translate"]

- **The problem of long sequences**

Ce genre d'expérience fait partie des efforts de Disney pour "prolonger la durée de vie de ses séries et créer de nouvelles relations avec des publics via des plateformes numériques de plus en plus importantes", a-t-il ajouté.

This kind of experience is part of Disney's efforts to "extend the lifetime of its series and build new relationships with audiences via digital platforms that are becoming ever more important," he added.



- **Encoder/Decoder :**

- $a^{(T_x)}$ is supposed to memorize the whole sentence then translate it ;
 - human way : translate each part of a sentence at a time

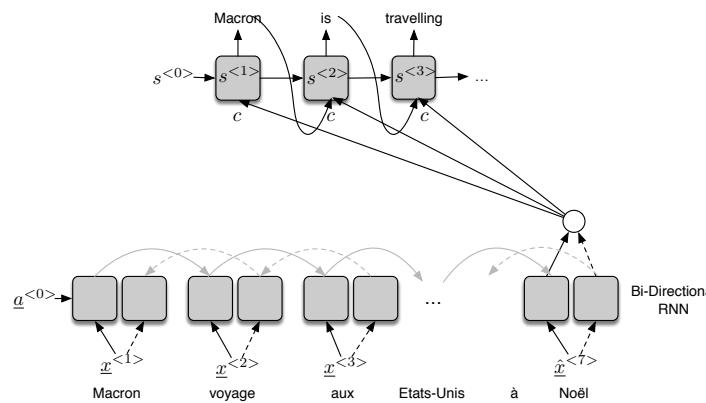
- **Attention model ?** (modification of the Encoder/Decoder)

- $a^{(T_x)}$ is replaced by a local version in the encoder \Rightarrow we pay attention on specific encoding

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 99

Sequence to sequence

Attention model

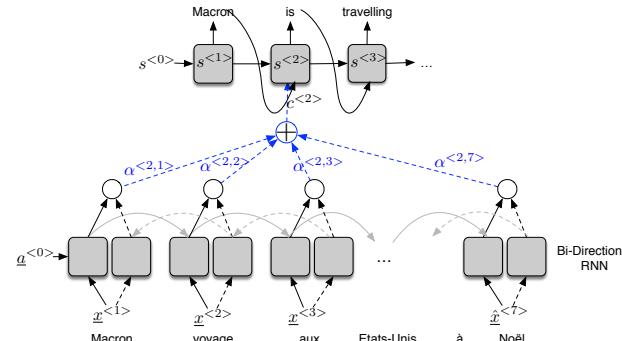


- In usual encoder-decoder,
 - information used for the decoding corresponds to the encoding at $T_x : c = a^{(T_x)}$

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 109

Sequence to sequence

Attention model

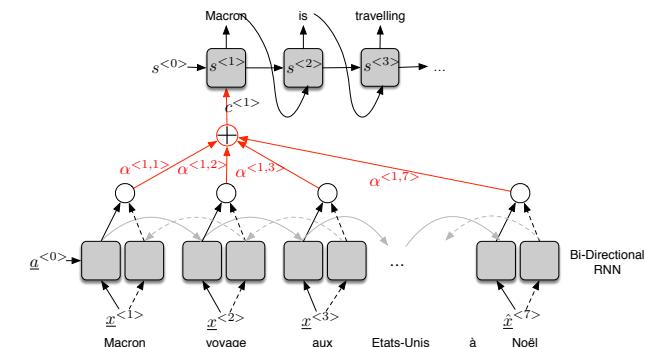


- Attention model
 - replace $c = a^{(T_x)}$ by a local version $c^{(t)}$
 - $c^{(t)}$ is computed as a weighted sum of the encoding hidden states $a^{(t)}$
- Attention weights $\alpha^{(\tau,t)}$:
 - when generating information at time $\tau = 1$, how much, do we have to pay attention on information at time $t = 1, 2, 3 \dots 7$
- Notation :
 - $a^{(t)} / s^{(\tau)}$: hidden states of encoder/ decoder

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 111

Sequence to sequence

Attention model

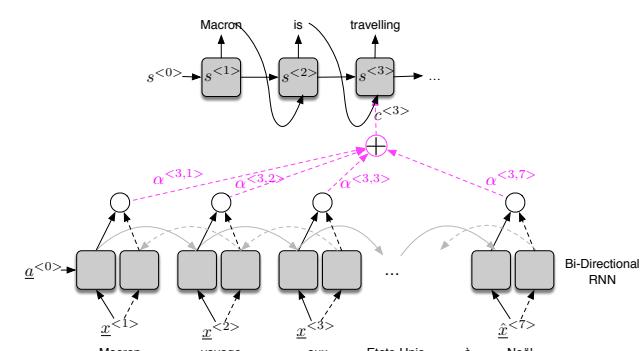


- Attention model
 - replace $c = a^{(T_x)}$ by a local version $c^{(t)}$
 - $c^{(t)}$ is computed as a weighted sum of the encoding hidden states $a^{(t)}$
- Attention weights $\alpha^{(\tau,t)}$:
 - when generating information at time $\tau = 1$, how much, do we have to pay attention on information at time $t = 1, 2, 3 \dots 7$
- Notation :
 - $a^{(t)} / s^{(\tau)}$: hidden states of encoder/ decoder

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 110

Sequence to sequence

Attention model



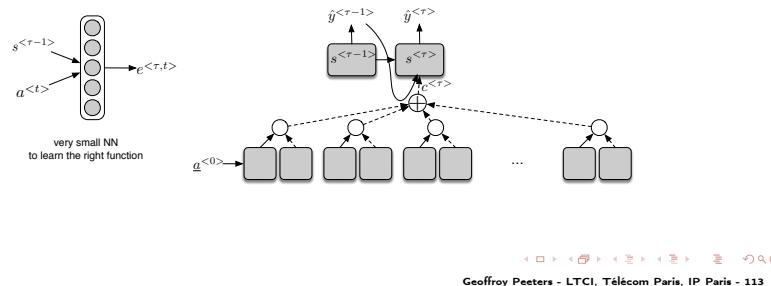
- Attention weights $\alpha^{(3,t)}$
 - describe an "alignment" between information at encoding time t and decoding time τ
 - encoding time t : computed using $\bar{a}^{(t)}$ and $\bar{a}^{(t)}$
 - we note $a^{(t)} = [\bar{a}^{(t)}, \bar{a}^{(t)}]$
 - decoding time $\tau = 3$: computed using $s^{(2)}$

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 112

Sequence to sequence

Attention model

- We define each conditional probability as
 - $p(y^{(\tau)}|y^{(1)}, \dots, y^{(\tau-1)}, x) = g(y^{(\tau-1)}, s^{(\tau)}, c^{(\tau)})$
 - where $s^{(\tau)}$ is an RNN hidden state for time τ
 - $s^{(\tau)} = f(s^{(\tau-1)}, y^{(\tau-1)}, c^{(\tau)})$
- For each target word $y^{(\tau)}$, the probability is conditioned on a distinct **context vector** $c^{(\tau)}$
 - $c^{(\tau)}$ depends on a sequence of annotations $(a^{(1)}, \dots, a^{(T_x)})$ to which an encoder maps the input sentence
 - Each $a^{(t)}$ contains information about the whole input sequence with a strong focus on the parts surrounding the t -th word of the input sequence



Sequence to sequence

Attention model

- Original sentence

This kind of experience is part of Disney's efforts to "extend the lifetime of its series and build new relationships with audiences via digital platforms that are becoming ever more important," he added.

- Encoder/Decoder without attention model

Ce type d'expérience fait partie des initiatives du Disney pour "prolonger la durée de vie de ses séries et créer de nouvelles relations avec des publics via des plateformes numériques de plus en plus importantes", a-t-il ajouté.

- Encoder/Decoder with attention model

Ce genre d'expérience fait partie des efforts de Disney pour "prolonger la durée de vie de ses séries et créer de nouvelles relations avec des publics via des plateformes numériques de plus en plus importantes", a-t-il ajouté.

Visualization of $\alpha^{(\tau,t)}$

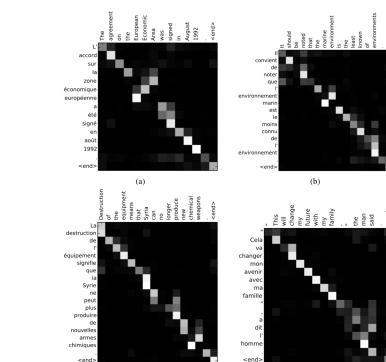
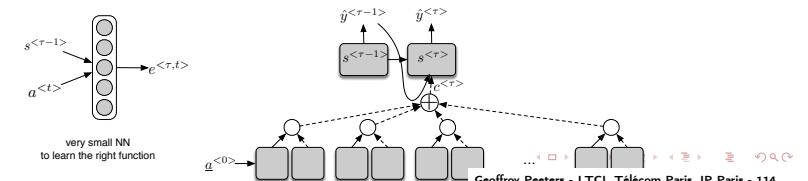


Figure 3: Four sample alignments found by RNN-seq-50. The x-axis and y-axis of each plot correspond to the words in the source sentence (English) and the generated translation (French), respectively. Each pixel shows the weight $\alpha_{ij}^{(\tau,t)}$ of the annotation of the i -th source word for the j -th target word (see Eq. (6)), in grayscale (0: black, 1: white). (a) an arbitrary sentence. (b-d) three random sentences from the test set. All four plots show sentences without any unknown words and of length between 10 and 20 words from the test set.

Sequence to sequence

Attention model

- The **context vector** $c^{(\tau)}$ at decoding time τ
 - computed as the sum of the annotations $a^{(t)}$ weighted by their **attention weights** $\alpha^{(\tau,t)}$
 - $c^{(\tau)} = \sum_t \alpha^{(\tau,t)} a^{(t)}$
- The **attention weight** $\alpha^{(\tau,t)}$
 - $\alpha^{(\tau,t)}$ = among of attention $y^{(\tau)}$ should pay to $a^{(t)}$
 - computed using a softmax on the **alignment model** $e^{(\tau,t)}$
 - $$\alpha^{(\tau,t)} = \frac{\exp(e^{(\tau,t)})}{\sum_{t'=1}^{T_x} \exp(e^{(\tau,t')})} \quad \text{with } \sum_t \alpha^{(\tau,t)} = 1$$
- The **alignment model** $e^{(\tau,t)}$
 - scores how well the inputs around position t match the output at position τ
 - computed using two inputs
 - the RNN hidden state $s^{(\tau-1)}$ (just before emitting $y^{(\tau)}$)
 - the t -th annotation $a^{(t)}$ of the input sentence.
 - computed using a feedforward neural network (jointly with all the other components)

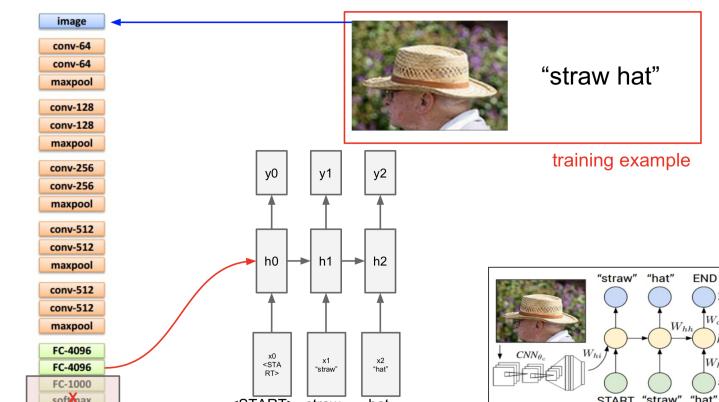


Sequence to sequence

Application : Image captioning

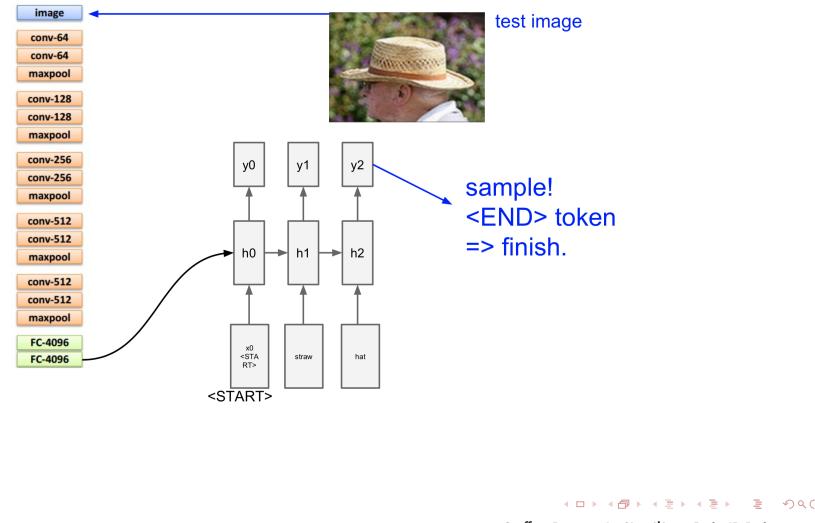
[Mao et al. 2014 "Deep captioning with multimodal recurrent neural networks (m-rnn)"
 [Vinyals et al. 2015 "Show and tell: A neural image caption generator"]
 [Karpathy et al. 2015 "Deep visual-semantic alignments for generating image descriptions"]

- <http://cs231n.stanford.edu>



Sequence to sequence

Application : Image captioning



Transformer

Sequence to sequence

Transformer

[Vaswani et al. 2017 "Attention Is All You Need"]

- <http://jamalmar.github.io/illustrated-transformer/>

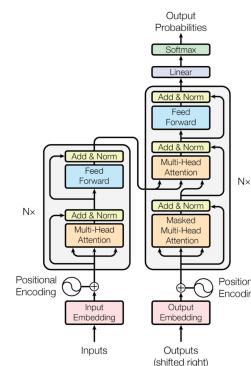
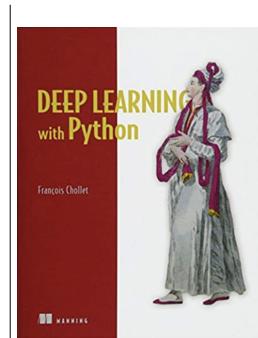


Figure 1: The Transformer - model architecture.

Implementations in keras

9- Implementations in keras 9.0-

- Readings :
 - <https://keras.io>
 - François Chollet : "Deep Learning with Python"
 - keras Cheat Sheet
 - Attilio Fiandratti keras slides



Dimensions : MLP

MLP in keras

```
[94] # --- Input dimensions are (m, n_in)
# --- Output dimensions are (m, n_out)
X_dummy = np.random.randn(20, 15, 100)
input_shape=[15, 100]

X_input = Input(shape = input_shape, name='Input')
X = Activation('relu', name='Activation1')(X_input)
X = Dense(1, name='FC2')(X)
X_output = Activation('sigmoid', name='Activation2')(X)

model = Model(inputs= X_input, outputs=X_output)

model.summary()
print('input dimension: ', X_dummy.shape)
print('output dimension: ', model.predict(X_dummy).shape)
```

Layer (type)	Output Shape	Param #
Input (InputLayer)	(None, 100)	0
FC1 (Dense)	(None, 128)	12928
Activation1 (Activation)	(None, 128)	0
FC2 (Dense)	(None, 1)	129
Activation2 (Activation)	(None, 1)	0

Total params: 13,057
Trainable params: 13,057
Non-trainable params: 0

input dimension: (20, 100)
output dimension: (20, 1)

Dimensions : RNN many-to-many

RNN in keras

```
[101] # --- Input dimensions are
# --- (m, seq_len1, n_in)
X_dummy = np.random.randn(20, 15, 100)
input_shape=[15, 100]

# --- (m, seq_len1) if use of embedding
#X_dummy = np.random.randint(0, 10000, (20,15))
#input_shape=[15,]

# --- Output dimensions are
# --- (m, seq_len1, n_out) if many-to-many Tx-Ty
# --- (m, seq_len2, n_out) if many-to-many Tx \ned Ty
# --- (m, n_out) if many-to-one

X_input = Input(shape = input_shape, name='Input')
#X = Embedding(10000, 300, name='embedding')(X_input) # --- 10000 is the size of the vocabulary / 300 is the size of the embedding
X, H, C = LSTM(32, return_sequences=True, return_state=True, name='RNN1')(X_input)
X, H, C = LSTM(32, return_sequences=True, return_state=True, name='RNN2')(X)
X = Dense(1, name='FC')(X)
X_output= Activation('sigmoid', name='Activation')(X)

model = Model(inputs= X_input, outputs=X_output)

model.summary()
print('input dimension: ', X_dummy.shape)
print('output dimension: ', model.predict(X_dummy).shape)
```

Layer (type)	Output Shape	Param #
Input (InputLayer)	(None, 15, 100)	0
RNN1 (LSTM)	[None, 15, 32], (None, 3 17024	
RNN2 (LSTM)	[None, 15, 32], (None, 3 8320	
FC (Dense)	(None, 15, 1)	33
Activation (Activation)	(None, 15, 1)	0

Total params: 25,377
Trainable params: 25,377
Non-trainable params: 0

input dimension: (20, 15, 100)
output dimension: (20, 15, 1)

Dimensions : RNN many-to-one

RNN in keras

```
[103] # --- Input dimensions are
# --- (m, seq_len1, n_in)
X_dummy = np.random.randn(20, 15, 100)
input_shape=[15, 100]

# --- (m, seq_len) if use of embedding
#X_dummy = np.random.randint(0, 10000, (20,15))
#input_shape=[15,]

# --- Output dimensions are
# --- (m, seq_len1, n_out) if many-to-many Tx-Ty
# --- (m, seq_len2, n_out) if many-to-many Tx \ned Ty
# --- (m, n_out) if many-to-one

X_input = Input(shape = input_shape, name='Input')
#X = Embedding(10000, 300, name='embedding')(X_input) # --- 10000 is the size of the vocabulary / 300 is the size of the embedding
X, H, C = LSTM(32, return_sequences=False, return_state=True, name='RNN1')(X_input)
X, H, C = LSTM(32, return_sequences=False, return_state=True, name='RNN2')(X)
X = Dense(1, name='FC')(X)
X_output= Activation('sigmoid', name='Activation')(X)

model = Model(inputs= X_input, outputs=X_output)

model.summary()
print('input dimension: ', X_dummy.shape)
print('output dimension: ', model.predict(X_dummy).shape)
```

Layer (type)	Output Shape	Param #
Input (InputLayer)	(None, 15, 100)	0
RNN1 (LSTM)	[None, 15, 32], (None, 3 17024	
RNN2 (LSTM)	[None, 32], (None, 32), 8320	
FC (Dense)	(None, 1)	33
Activation (Activation)	(None, 1)	0

Total params: 25,377
Trainable params: 25,377
Non-trainable params: 0

input dimension: (20, 15, 100)
output dimension: (20, 1)

Dimensions : RNN many-to-one with embedding

RNN in keras

```
[104] # --- Input dimensions are
# --- (n_in, n_out, seq_len, n_in)
# x_dummy = np.random.randint(20, 15, 100)
#input_shape=(15, 100,)

# --- (m, seq_len) if use of embedding
x_dummy = np.random.randint(0, 10000, (20,15))
input_shape=(15,)

# --- Output dimensions are
# --- (m, seq_len, n_out) if many-to-many TxTy
# --- (m, seq_len, n_out) if many-to-many Tx \neq Ty
# --- (m, n_out) if many-to-one

x_input = Input(shape = input_shape, name='Input')
# --- Embedding(10000, 300, name='Embedding')(x_input) # --- 10000 is the size of the vocabulary / 300 is the size of the embedding
x, B, C = LSTM(32, return_sequences=True, return_state=True, name='RNN1')(x)
x, B, C = LSTM(32, return_sequences=False, return_state=True, name='RNN2')(x)
x_output= Activation('sigmoid', name='Activation')(x)

model = Model(inputs= x_input, outputs=x_output)

model.summary()
print('input dimension: ', x_dummy.shape)
print('output dimension: ', model.predict(x_dummy).shape)
```

Layer (type)	Output Shape	Param #
Input (InputLayer)	(None, 15)	0
Embedding (Embedding)	(None, 15, 300)	3000000
RNN1 (LSTM)	[None, 15, 32], (None, 3 42624	
RNN2 (LSTM)	[None, 32], (None, 32), 8320	
FC (Dense)	(None, 1)	33
Activation (Activation)	(None, 1)	0

Total params: 3,650,977
Trainable params: 3,650,977
Non-trainable params: 0

input dimension: (20, 15)
output dimension: (20, 1)

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 142

Dimensions : RNN padding sequences of different lengths

Keras Padding Sequences

```
from keras.preprocessing.sequence import pad_sequences
from keras.preprocessing.text import Tokenizer

reviews = ['This movie is fantastic ! I really like it because it is so good !',
           'Not to my taste, will skip and watch another movie.',
           'This movie was completely lacking a good script, good actors and a good director.']

data = [review.split() for review in reviews]
print('*** Sequences have different lengths: ', print(len(data[0])), print(len(data[1])), print(len(data[2])))

# --- First tokenize all the texts
tokenizer_obj = Tokenizer()
tokenizer_obj.fit_on_texts(reviews)
print('*** Tokenized texts: ')
print(tokenizer_obj.word_index)
print('*** Index to word: ', tokenizer_obj.index_word)

x_train_tokens = tokenizer_obj.texts_to_sequences(reviews)
print('*** Tokenized texts: ')
print(x_train_tokens)
print('*** Sequences have different lengths')

max_length = 40
X_train_pad = pad_sequences(x_train_tokens, maxlen=max_length, value=0.0, padding='post')
print('*** Post-padding: ')
print(X_train_pad)
X_train_pad = pad_sequences(x_train_tokens, maxlen=max_length, value=0.0, padding='pre')
print('*** Pre-padding: ')
print(X_train_pad)

Geoffroy Peeters - LTCI, Télécom Paris, IP Paris - 143
```