

# TP2: Word Embedding - RNN

Hachem Madmoun

November 4, 2019



# OUTLINE

## 1 Recurrent Neural Networks - LSTM

- Introduction
- The idea behind Recurrent Neural Networks
- LSTM
- Step-by-Step LSTM Walk Through

## 2 LSTM - Applications

- Categories of tasks of LSTM Model
- Many to one : Spam Detection

## 3 Iteration Based Methods - Word2vec

- Introduction
- Main ideas of word2vec
- Steps of Skip Gram



# OUTLINE

## 1 Recurrent Neural Networks - LSTM

- Introduction
- The idea behind Recurrent Neural Networks
- LSTM
- Step-by-Step LSTM Walk Through

## 2 LSTM - Applications

- Categories of tasks of LSTM Model
- Many to one : Spam Detection

## 3 Iteration Based Methods - Word2vec

- Introduction
- Main ideas of word2vec
- Steps of Skip Gram

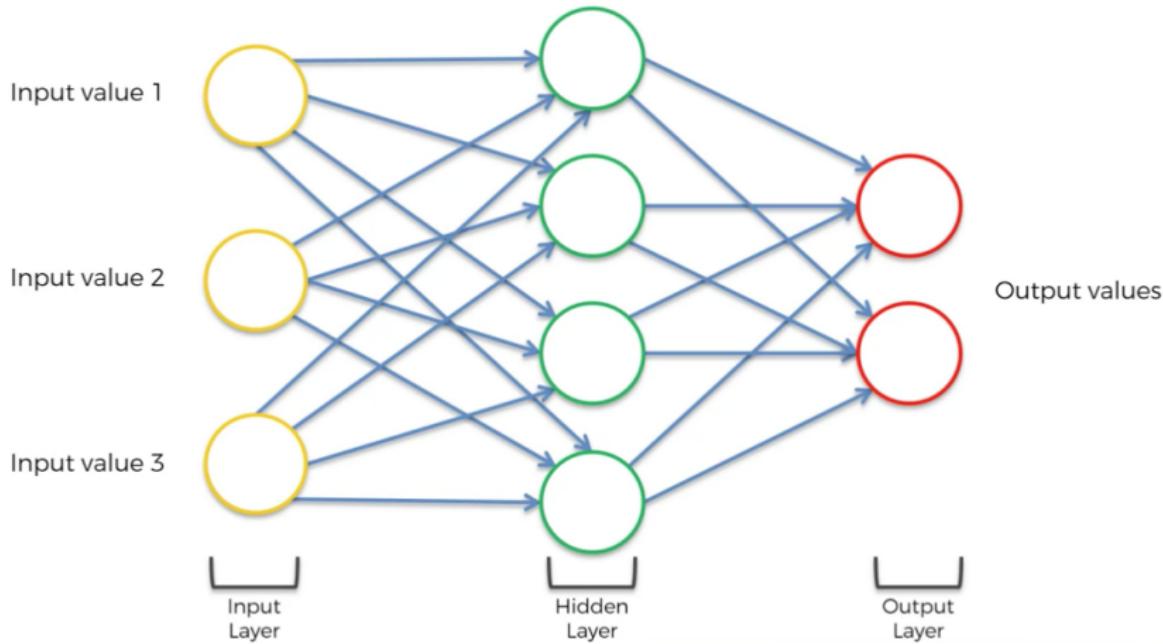


# PLAN OF ATTACK

- ▶ The idea behind Recurrent Neural Networks
- ▶ The vanishing Gradient Problem
- ▶ LSTM
- ▶ Step-by-Step LSTM Walk Through



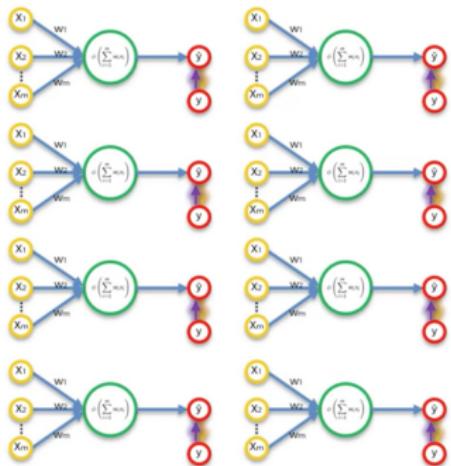
# A NORMAL NEURAL NETWORK



○  
○○○○

10

# HOW DO NEURAL NETWORKS LEARN?



Row ID	Study Hrs	Sleep Hrs	Quiz	Exam
1	12	6	78%	93%
2	22	6.5	24%	68%
3	115	4	100%	95%
4	31	9	67%	75%
5	0	10	58%	51%
6	5	8	78%	60%
7	92	6	82%	89%
8	57	8	91%	97%

$$C = \sum \frac{1}{2}(\hat{y} - y)^2$$



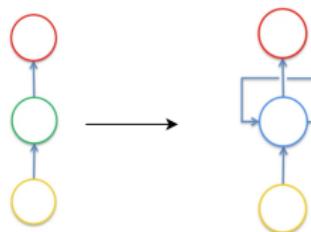
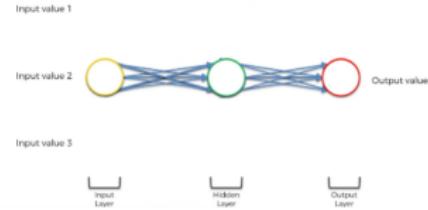
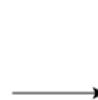
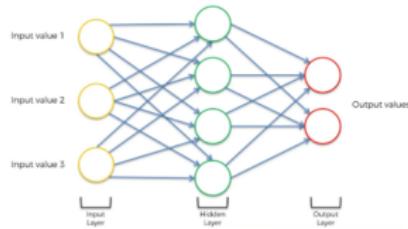


## BACKPROPAGATION

- ▶ **Step 1 :** Randomly initialise the weights to small numbers close to 0 (NOT 0)
- ▶ **Step 2 :** Input the first observation in the input layer, each feature in one input node.
- ▶ **Step 3 :** Forward-Propagation : from left to right, the neurons are activated in a way that the impact of each neuron's activation is limited by the weights. Propagate the activations until getting the predicted result
- ▶ **Step 4 :** Back-propagation : from right to left, the error is back-propagated. Update the weights, the learning rate decides by how much we update the weights.
- ▶ **Step 5 :** Repeat Steps 1 to 4 and update the weights after each observation  $\Rightarrow$  Reinforcement Learning / after a batch of observation  $\Rightarrow$  Batch Learning.
- ▶ **Step 6 :** When the whole training set is passed through the ANN, that makes an epoch. Redo more epochs

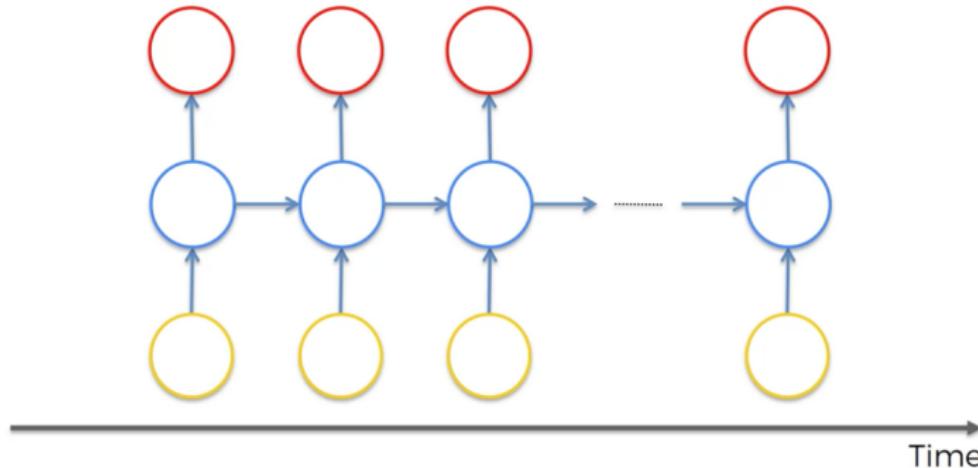


# FROM ANN TO RNN



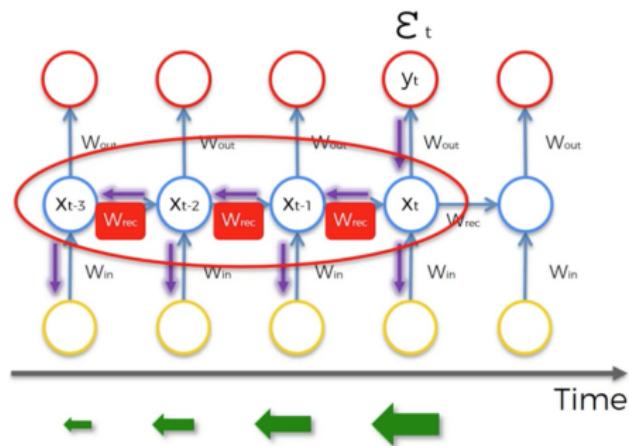


# RNN





# THE VANISHING GRADIENT PROBLEM



$$\frac{\partial \mathcal{E}}{\partial \theta} = \sum_{1 \leq t \leq T} \frac{\partial \mathcal{E}_t}{\partial \theta} \quad (3)$$

$$\frac{\partial \mathcal{E}_t}{\partial \theta} = \sum_{1 \leq k \leq t} \left( \frac{\partial \mathcal{E}_t}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} \frac{\partial \mathbf{x}_k}{\partial \theta} \right) \quad (4)$$

$$\frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \prod_{i \geq i > k} \frac{\partial \mathbf{x}_i}{\partial \mathbf{x}_{i-1}} = \prod_{i \geq i > k} \mathbf{W}_{rec}^T diag(\sigma'(\mathbf{x}_{i-1})) \quad (5)$$

$W_{rec} \sim \text{small}$	→	Vanishing
$W_{rec} \sim \text{large}$	→	Exploding

Formula Source: Razvan Pascanu et al. (2013)

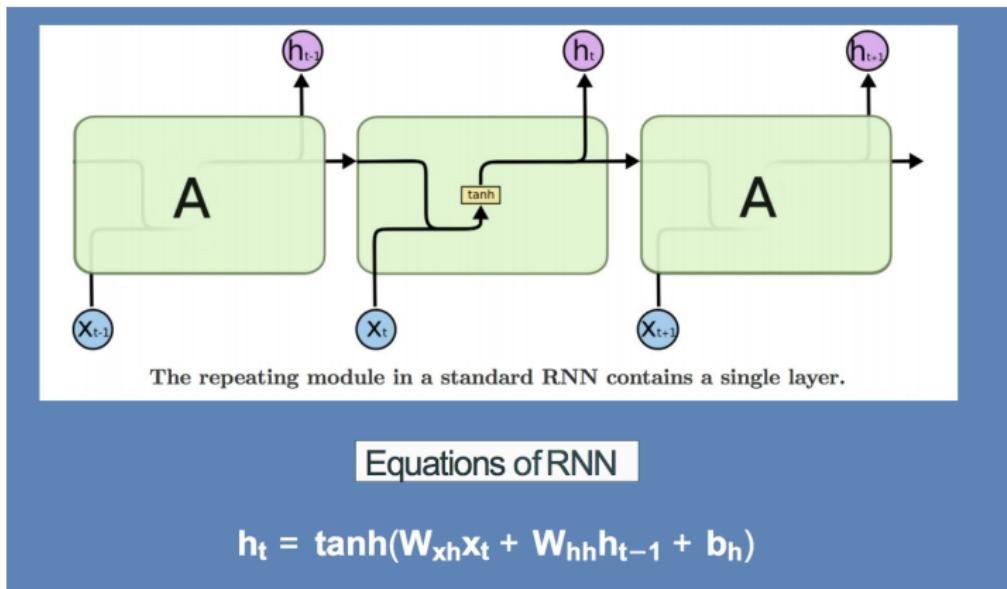


# RESOURCES

- ▶ Untersuchungen zu dynamischen neuronalen Netzen (1991) by Sepp (Josef) Hochreiter  
[http://people.idsia.ch/~juergen/  
SeppHochreiter1991ThesisAdvisorSchmidhuber.pdf](http://people.idsia.ch/~juergen/SeppHochreiter1991ThesisAdvisorSchmidhuber.pdf)
- ▶ Learning Long-Term Dependencies with Gradient Descent is Difficult (1994) by Yoshua Bengio.  
[http://www.iro.umontreal.ca/~lisa/poiteurs/  
ieeetrnn94.pdf](http://www.iro.umontreal.ca/~lisa/poiteurs/ieeetrnn94.pdf)
- ▶ On the difficulty of training recurrent neural networks (2013) by Razvan Pascanu and Yoshua Bengio.  
<http://proceedings.mlr.press/v28/pascanu13.pdf>



# STANDARD RNN





# LSTM

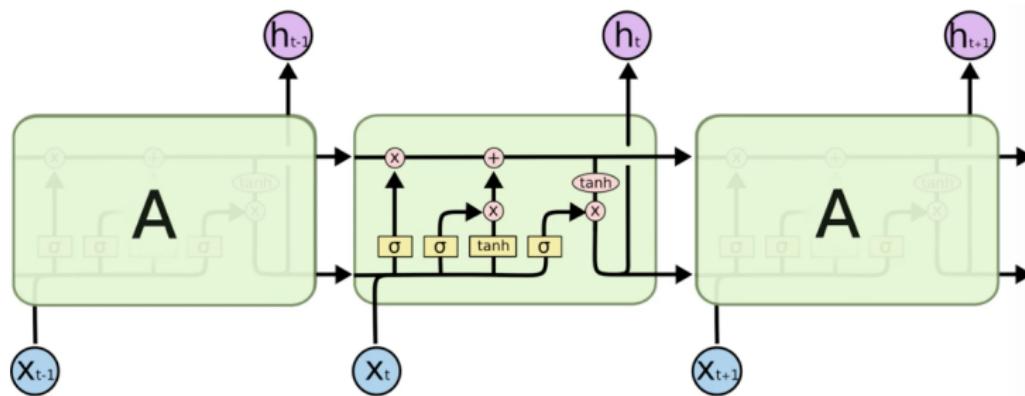
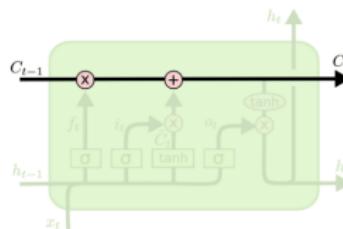


Image Source: colah.github.io

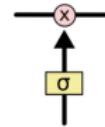


# THE CORE IDEA BEHIND LSTM

- ▶ The cell state :



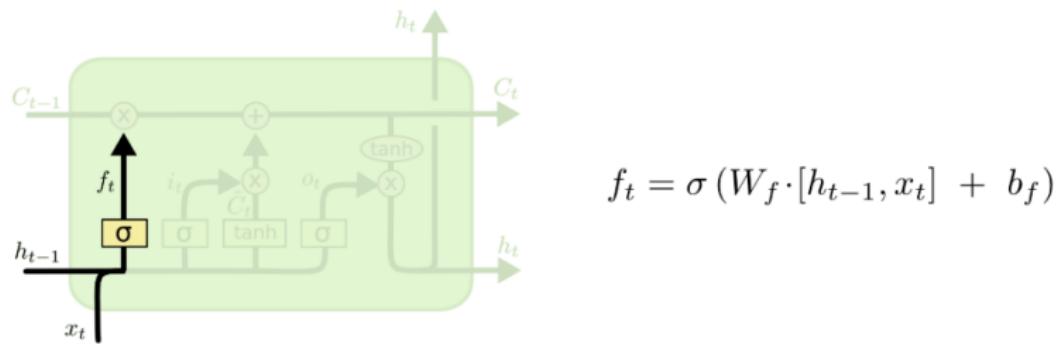
- ▶ The concept of gate :



- ▶ The LSTM has three of these gates.

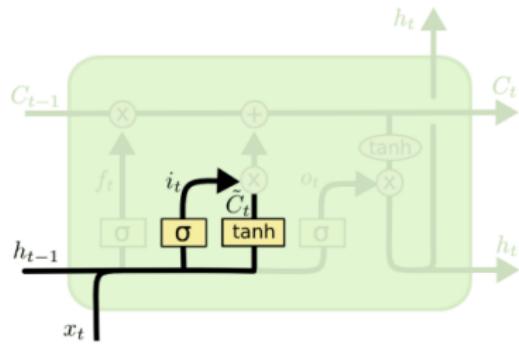


## STEP 1 : THE FORGET GATE LAYER





## STEP 2 : THE INPUT GATE LAYER



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

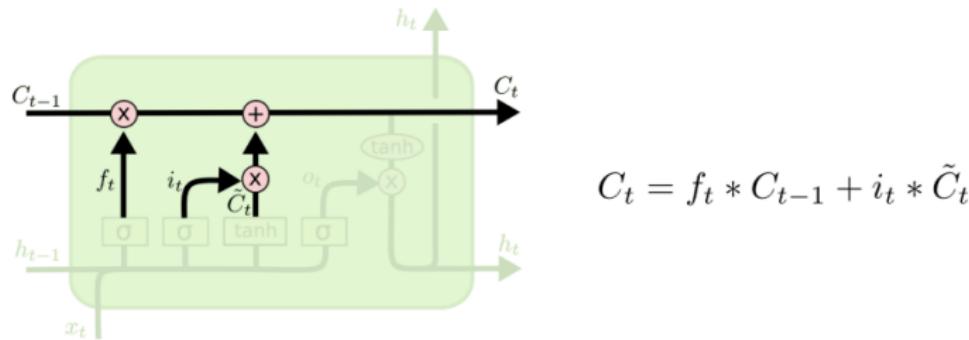
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

○  
 ○○○○○○○  
 ○○○○  
 ○○●○○

○  
 ○○○○

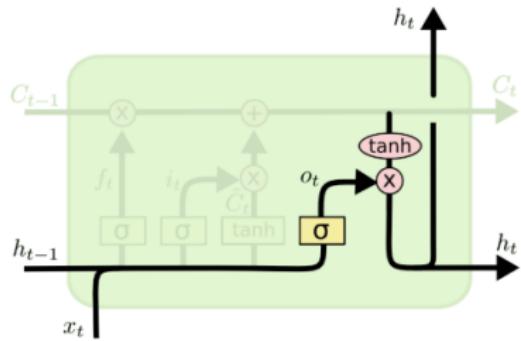
○  
 ○○○○  
 ○○○○○

## STEP 3 : UPDATE THE CELL STATE





## STEP 4 : THE OUTPUT GATE

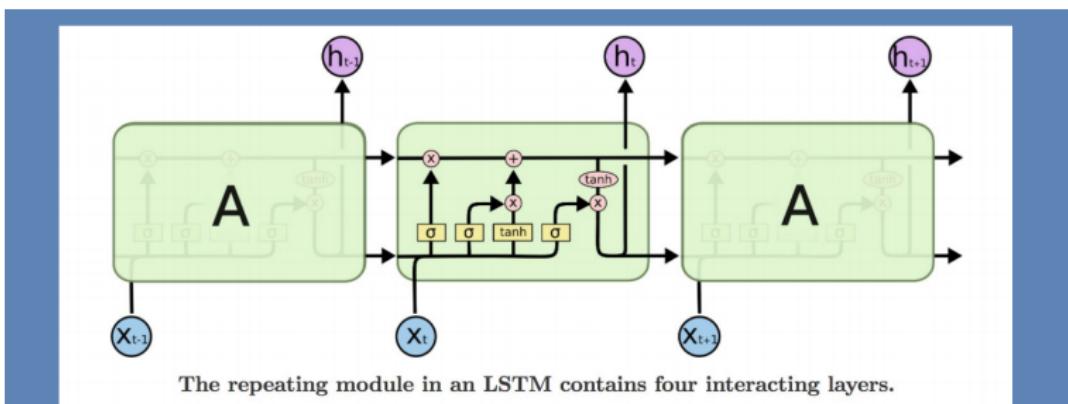


$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$



# SUMMARY LSTM



## Equations of LSTM

### Gates

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i)$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + b_f)$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + b_o)$$

### Updates

$$\tilde{c}_t = \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c)$$

$$c_t = f_t c_{t-1} + i_t \tilde{c}_t$$

$$h_t = o_t \tanh(c_t)$$



# OUTLINE

## 1 Recurrent Neural Networks - LSTM

- Introduction
- The idea behind Recurrent Neural Networks
- LSTM
- Step-by-Step LSTM Walk Through

## 2 LSTM - Applications

- Categories of tasks of LSTM Model
- Many to one : Spam Detection

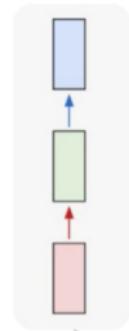
## 3 Iteration Based Methods - Word2vec

- Introduction
- Main ideas of word2vec
- Steps of Skip Gram

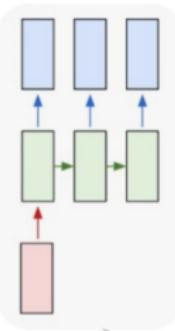


# CATEGORIES OF TASKS

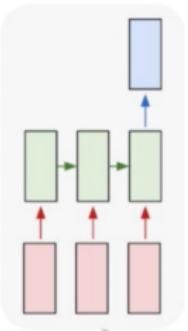
one to one



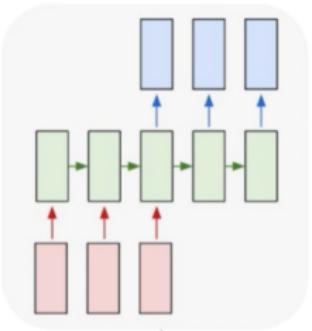
one to many



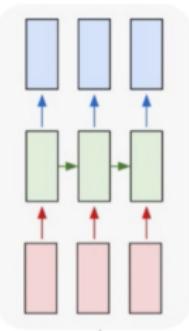
many to one



many to many



many to many



Feedforward

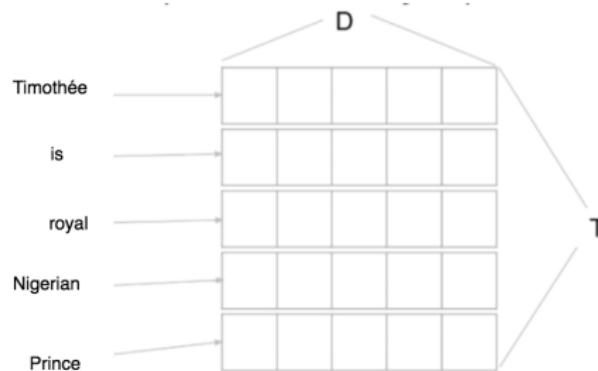
Poetry generation

Spam detection  
Sentiment analysisMachine Translation  
Chatbots  
Question AnsweringParts of speech  
Named entity rec.



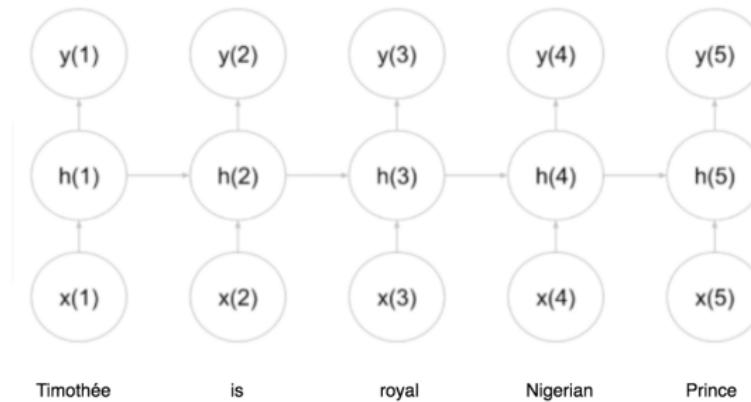
## RNN INPUT

- ▶ The input is a sentence, each token is encoded in  $D$ -dimensional vector.
- ▶  $T$  : is the sequence length,  $D$  : input dimensionality.



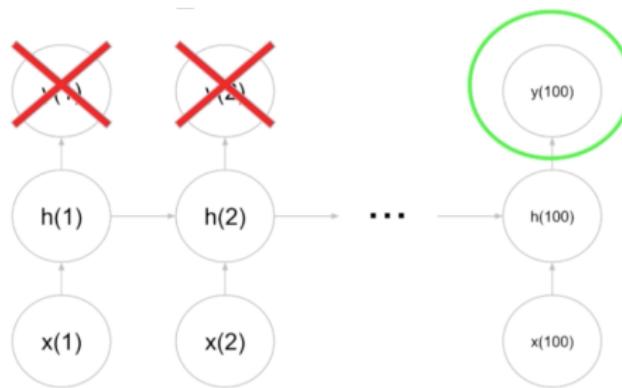


# RNN OUTPUT





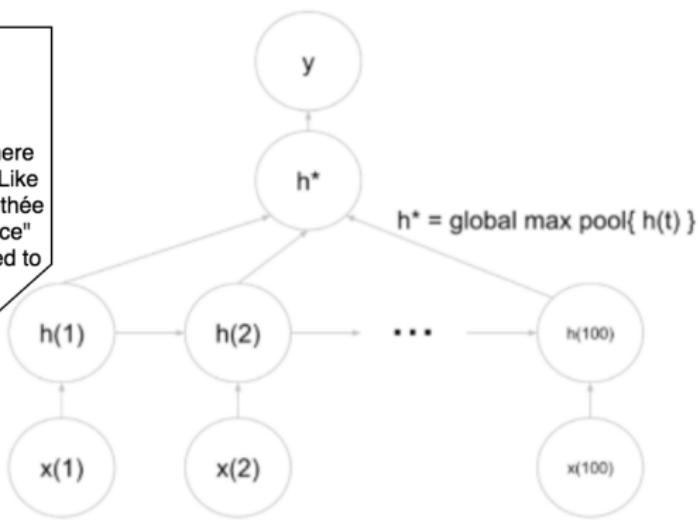
## FIRST CASE - WE TAKE THE LAST OUTPUT





## SECOND CASE - WE USE A GLOBAL MAXPOOL

Why ? We don't know where the useful state may be. Like when it contains "I'm Timothée ! I'm a royal nigerian prince" It's clearly a spam, no need to read the whole text.





# OUTLINE

## 1 Recurrent Neural Networks - LSTM

- Introduction
- The idea behind Recurrent Neural Networks
- LSTM
- Step-by-Step LSTM Walk Through

## 2 LSTM - Applications

- Categories of tasks of LSTM Model
- Many to one : Spam Detection

## 3 Iteration Based Methods - Word2vec

- Introduction
- Main ideas of word2vec
- Steps of Skip Gram



## SKIP-GRAM - CBOW

- ▶ The idea behind this method is that we can get a lot of value by representing a word by means of its neighbors. As John Rupert Firth<sup>1</sup> said : "*You shall know a word by company it keeps*"
- ▶ So, we will build a dense vector for each word type, chosen so that it is good at predicting other words appearing in its context.
- ▶ There are two Algorithms :
  - ➊ Skip-Gram (SG) : Predict context words given target (position independent)
  - ➋ Continuous Bag of Words (CBOW) : Predict Target word from bag-of-words context

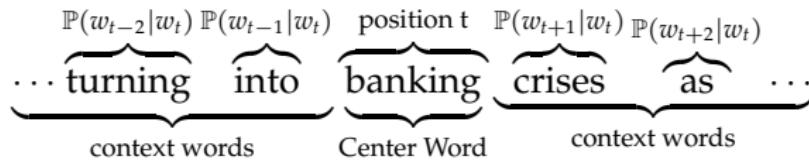
---

1. Leading figure in British linguistics during 1950s



# WORD2VEC

- ▶ The aim is to create a model such that given the center word, the model will be able to predict or generate the surrounding words.
- ▶ For example :



- ▶ So, for each word  $t = 1, \dots, T$ , we want to predict the surrounding words in a window of radius m.



# WORD2VEC

- ▶ **Each word has two vector representations<sup>2</sup>**, one as an input (which is the center word for the skip gram model), and an other as an output (which is the context for this model) : " $u_w$ " will represent the vector form of the word "w" as a context word and " $v_w$ " will represent its vector form as a center word.
- ▶ So, there are two matrices that we should learn :  $\mathcal{V} \in \mathbb{R}^{N \times |V|}$  the **Input Word Matrix** and  $\mathcal{U} \in \mathbb{R}^{|V| \times N}$  the **Output Word Matrix**.
- ▶ We should notice that there are some hyperparameters :  $m$  the window size and the dimension of the vectors.

- 
2. We could have one word vector representation, but by using two word vector representations, it turns out that not only does that make the math easier because the two representations are separated when performing optimization, but it actually works better in practise



# THE OPTIMIZATION PROBLEM

- ▶ Let  $\theta$  represent all the vector representations of the words, which we aim to optimize.
- ▶ We will see that the optimization problem is :

$$\text{Maximize} \quad \prod_{t=1}^T \prod_{\substack{j \neq 0 \\ -m \leq j \leq m}} \mathbb{P}(w_{t+j}|w_t; \theta)$$

- ▶ By taking the logarithm. And rather than having the probability of the whole corpus, we can take the average over each position<sup>3</sup>. Thus, the optimization problem becomes :

$$\text{Minimize} \quad -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{j \neq 0 \\ -m \leq j \leq m}} \log(\mathbb{P}(w_{t+j}|w_t; \theta))$$

- 
3. The normalization won't affect the maximum



# THE OPTIMIZATION PROBLEM

For  $\mathbb{P}(w_{t+j}|w_t; \theta)$ , we will use the softmax form :

- ▶ Let "o" be the output word index and "c" the center word index. So, " $v_c$ " and " $u_o$ " are "center" and "output" vectors of indices c and o

$$\mathbb{P}(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^{|V|} \exp(u_w^T v_c)}$$

- ▶ The dot product  $u^T v$  gets bigger if u and v are "more similar".
- ▶ the softmax form is the standard way to turn numbers into probabilities (exponentiate them to make them positive and normalize them to give probability).

# THE STEPS OF SKIP GRAM MODEL

- ▶ We first generate our one-hot vector  $x \in \mathbb{R}^{|N|}$
- ▶ We get our **embedded vector**  $v_c = \mathcal{V}x \in \mathbb{R}^N$

$$N \begin{pmatrix} |V| & & & & |V| \\ v_1 & \cdots & v_c & \cdots & v_{|V|} \\ | & | & | & | & | \\ v_1 & \cdots & v_c & \cdots & v_{|V|} \end{pmatrix} \begin{pmatrix} |V| \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} = v_c$$



## THE STEPS OF SKIP GRAM MODEL

- We generate our score vector  $z = \mathcal{U}v_c$  with :  $y_o = u_o^T v_c$

$$|V| \begin{pmatrix} N \\ - & u_1 & - \\ \vdots \\ - & u_o & - \\ \vdots \\ - & u_{|V|} & - \end{pmatrix} v_c = \begin{pmatrix} |V| \\ y_1 \\ \vdots \\ y_o \\ \vdots \\ y_{|V|} \end{pmatrix}$$

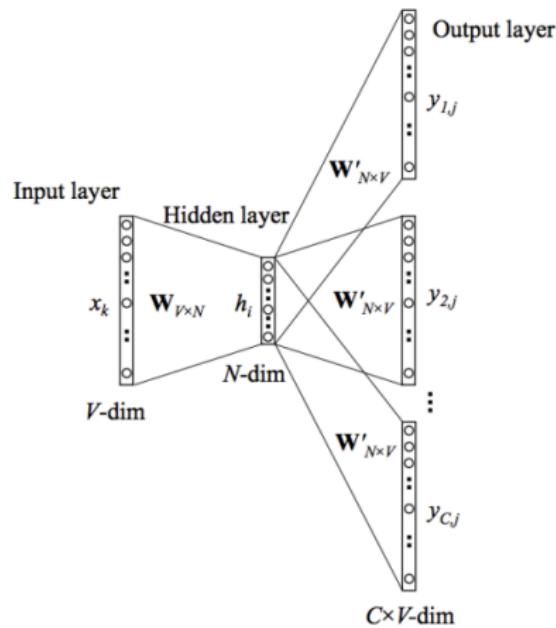
- We turn the vector into probabilities  $\hat{y} = \text{softmax}(z)$  :

$$\hat{y}_o = \mathbb{P}(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^{|V|} \exp(u_w^T v_c)}$$



# SKIP GRAM - DIAGRAM

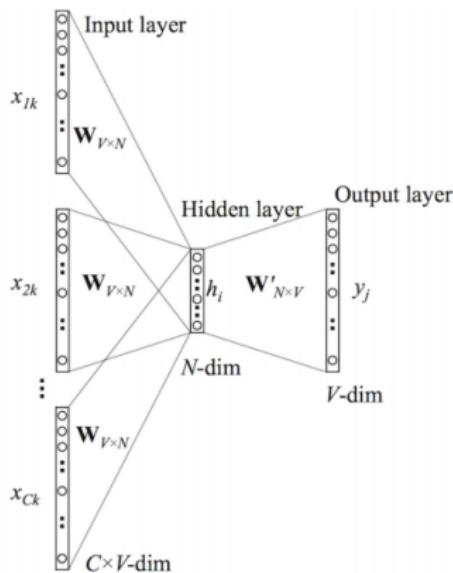
Below is the diagram of the skip gram model :





## CBOW - DIAGRAM

Below is the diagram of the CBOW model, which is like the Skip-Gram Model with the input and output reversed. So, the objective is to predict a center word from the surrounding context.



○  
○○○○○○○  
○○○  
○○○○

○  
○○○○

○  
○○○○  
○○○●

Thanks for your attention