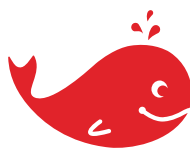


Formation

# Introduction Deep Learning

Séquence 05b

Données creuses / textuelles (Embedding)



FIDLE

<https://fidle.cnrs.fr>

5b



**Sparse data (text)**  
Embedding

5.1

## Text encoding

- From text to tensor
- One hot encoding
- Embedding

5.2

## Example 1 : IMDB1

- Sentiment analysis with one-hot encoding



5.3

## Example 2 : IMDB2/3/4

- Sentiment analysis with embedding



5b



**Sparse data (text)**  
Embedding

5.1

## Text encoding

- From text to tensor
- One hot encoding
- Embedding

!

Please note that this is an extremely large subject and this is only an introduction !

And mostly outdated, because progress is very fast ;-)

5.2

## Example 1 : IMDB1

- Sentiment analysis with one-hot encoding



5.3

## Example 2 : IMDB2/3/4

- Sentiment analysis with embedding

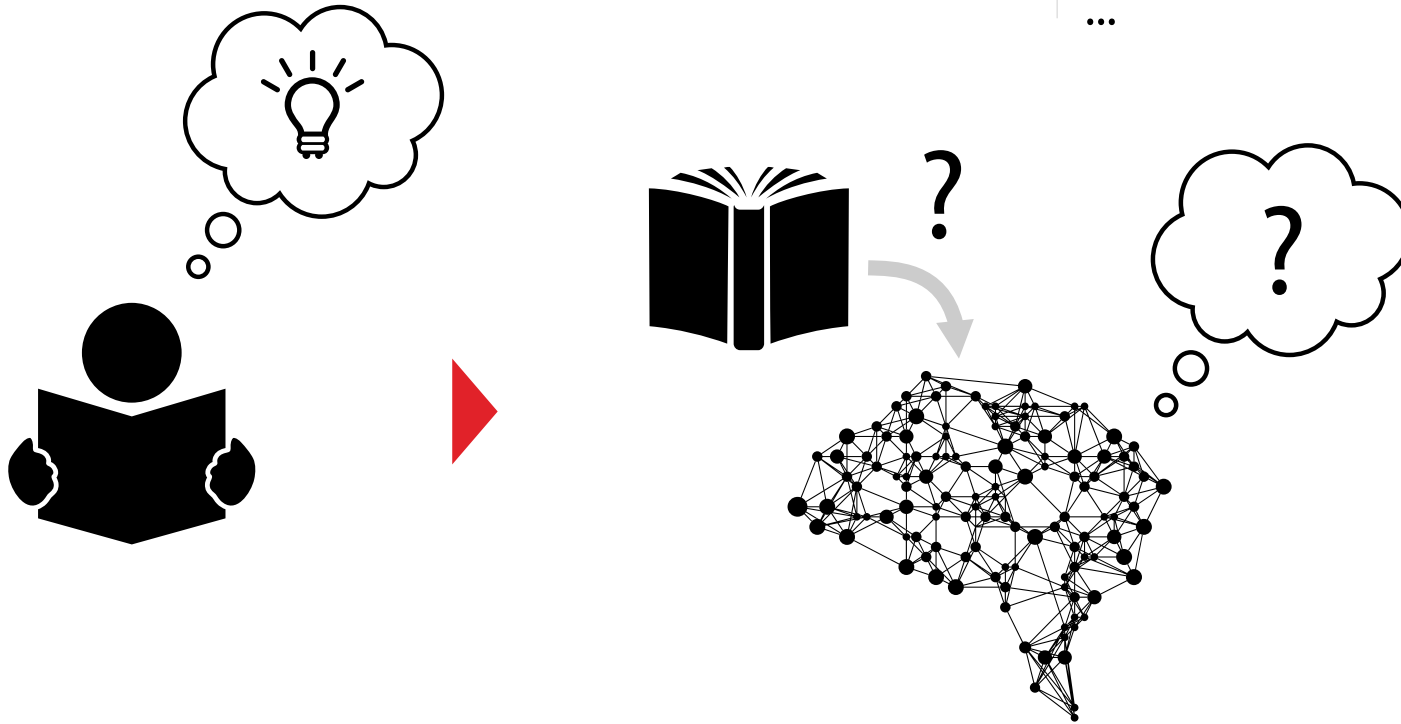


# How to feed a neural network with text?



Variable size  
Large size space  
Conceptual contents ;-)

...



« I've never seen a movie like this before. »



How to build a descriptor for this kind of data (text, DNA, ...)?

« I've never seen a movie like this before. »

Dictionary

0	a
1	before
2	fantastic
3	i've
4	is
5	like
6	movie
7	never
8	seen
9	this

["i've", 'never', 'seen', 'a', 'movie', 'like', 'this', 'before']



The word order in the dictionary may have some meaning. For example according to their rate of use.



Tokenization  
(words)

« I've never seen a movie like this before. »

Dictionary

0	a
1	before
2	fantastic
3	i've
4	is
5	like
6	movie
7	never
8	seen
9	this

["i've", 'never', 'seen', 'a', 'movie', 'like', 'this', 'before']

[ 3, 7, 8, 0, 6, 5, 9, 1 ]



The values associated with each word are **meaningless**: words with a contiguous subscript are typically unrelated.



Tokenization  
(words)



Vectorization

# One hot encoding

« I've never seen a movie like this before. »

["i've", 'never', 'seen', 'a', 'movie', 'like', 'this', 'before']

Dictionary

0	a
1	before
2	fantastic
3	i've
4	is
5	like
6	movie
7	never
8	seen
9	this

[ 3, 7, 8, 0, 6, 5, 9, 1 ]

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



Tokenization  
(words)



Vectorization

001000

One hot  
encoding

Each vector is independent. Each word has its own dimension.



# One hot encoding

« I've never seen a movie like this before. »

["i've", 'never', 'seen', 'a', 'movie', 'like', 'this', 'before']

Dictionary

0	a
1	before
2	fantastic
3	i've
4	is
5	like
6	movie
7	never
8	seen
9	this

[ 3, 7, 8, 0, 6, 5, 9, 1 ]

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



Each vector is independent. Each word has its own dimension.

# One hot encoding

## Limits of full one-hot encoding : Size !



Example :

Dictionary = 80 000 words

Sentence = 300 words



0	0	0	1	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	1
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	0
0	0	0	0	1	0	0	0
0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0

...

0	0	1	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0
0	0	0	0	0	0	0	0



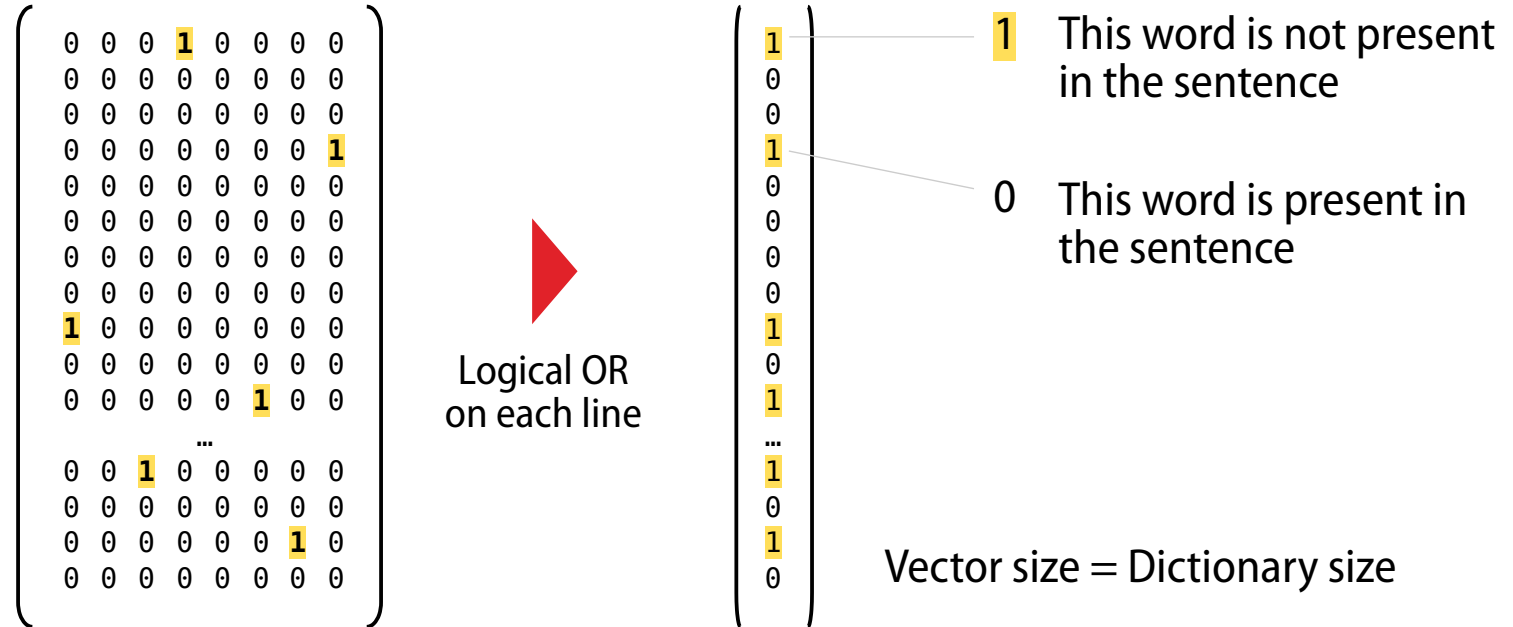
2 answers: { One hot vector  
Embedding

Need :  
24  $10^6$   
Parameters !



# One hot vector encoding

## Solution 1 : one hot vector



The size no longer depends on the length of the sentence.  
The size is drastically reduced, but we lose the words order.

## Solution 2 : Using embedding $\rightleftharpoons$

(Word<sub>i</sub>)  $\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \dots \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

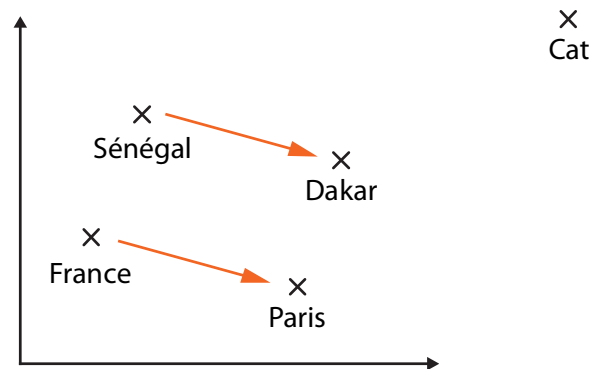


(W<sub>i</sub>)  $\begin{pmatrix} 0.7 \\ 4.2 \\ 1.6 \end{pmatrix}$

Short dense vector, whose value carries a semantic meaning of the word.

Long sparse vector, whose value just indicates membership in the dictionary

Semantic meaning, allows to perform calculations, such as :



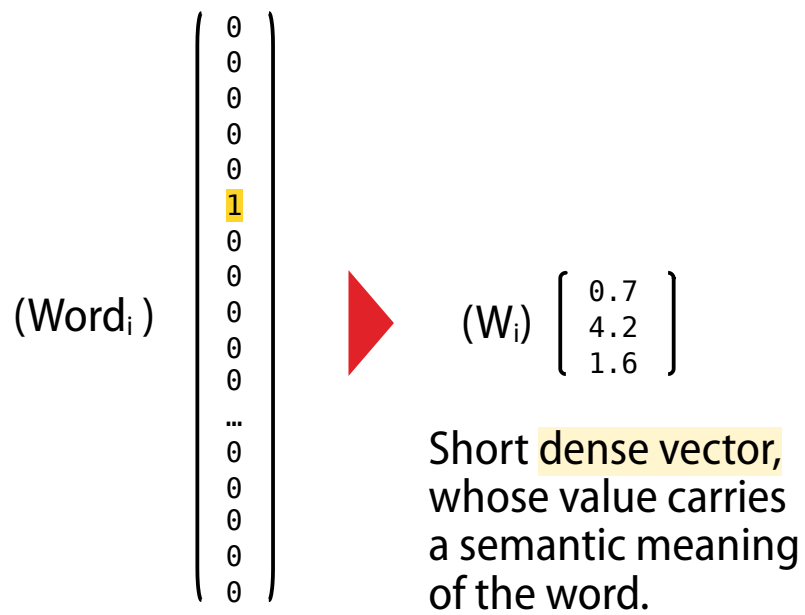
$\text{dist}(\text{"France"}, \text{"Cat"}) \gg \text{dist}(\text{"France"}, \text{"Sénégal"})$

$\text{vect}(\text{"France"}, \text{"Paris"}) \approx \text{vect}(\text{"Sénégal"}, \text{"Dakar"})$

$\text{"Paris"} = \text{"France"} + \text{vect}(\text{"Sénégal"}, \text{"Dakar"})$

*« Paris est à la France ce que Dakar est au Sénégal »*

## Solution 2 : Using embedding




Long sparse vector, whose value just indicates membership in the dictionary


Number of parameters required :

Example :

Embedding size : 200

Sentence size : 300 words 


Need only : 60.000 Parameters :-)



Embedding techniques :

- Contextual embedding { Keras embedding
- Global embedding { CBOW, SG, GloVe, etc.

## Embedding layers in Keras



Usable as a **simple layer**

This layer will constitute a **dictionary of vectors** which will be **optimized during the learning process**, according to the expected result and not to the pure semantics.

Keras embedding is therefore adapted to **classification**, but will not be able to take into consideration, for example, semantic similarities (such as identifying 2 sentences with the same semantic meaning).

The **output** of the layer is a **set of vectors**.

## Embedding layers in Keras

(french)



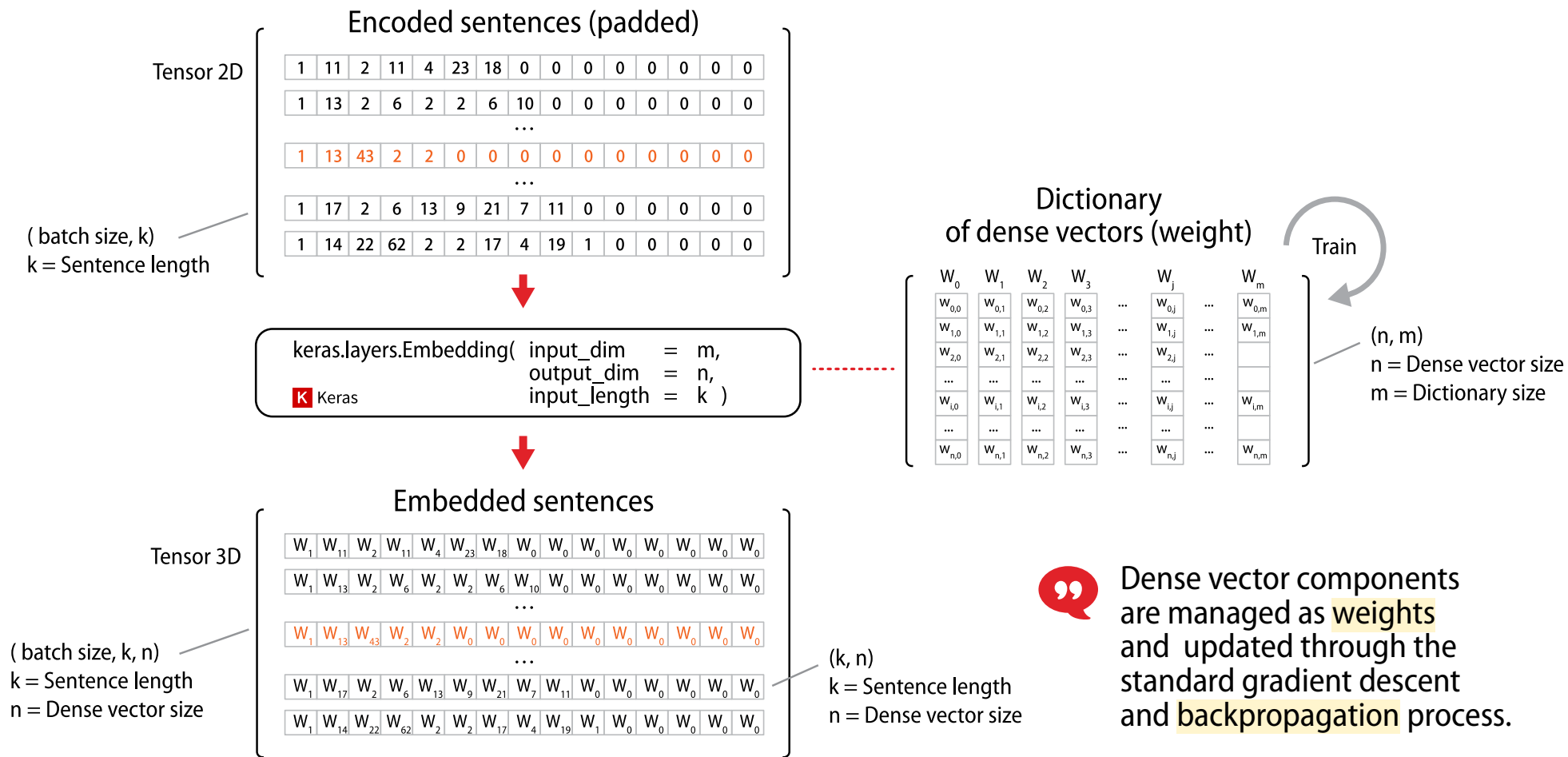
Utilisable comme une **simple couche**

Cette couche va constituer un **dictionnaire de vecteurs** qu'elle optimisera au cours de l'apprentissage, en **fonction du résultat** attendu et non de la sémantique pure.

L'embedding Keras est donc adapté à la **classification**, mais ne rendra pas compte, par exemple, de similarités sémantiques (comme identifier 2 phrases ayant un même sens sémantique).

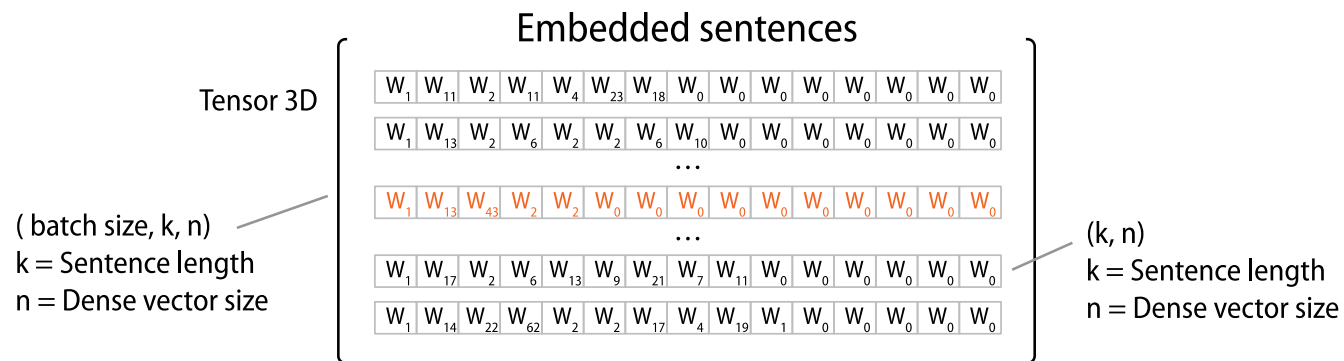
La **sortie** de la couche est un **ensemble de vecteurs**

# Keras embedding layer





# Keras embedding layer



`keras.layers.GlobalAveragePooling1D()`

 Keras

$$\overline{W^{(s)}} = \frac{1}{k} \sum_{i=1}^k W_i^{(s)}$$

Where :

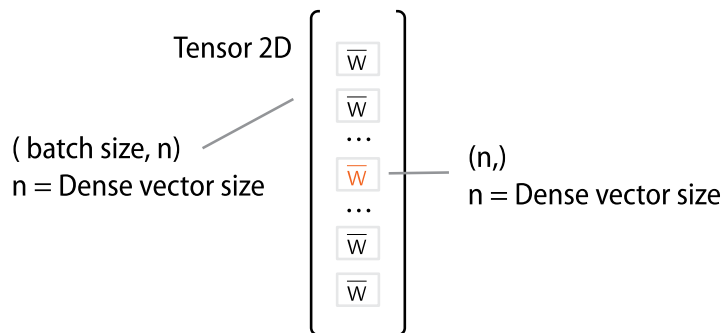
$k$  is sentence length

$W_i^{(s)}$  is word vector  $i$  of sentence  $(s)$

Each sentence is here encoded by a **unique dense vector** which is the average of the words composing the sentence.



Dense vector sentences



The **order** of the words is therefore **ignored** !

## Word2Vec<sup>1</sup>

This approach aims to build **dictionaries** whose vector representation of words is based on **context** and therefore on **semantics**.

Dictionaries built from large corpora are available.

Two models:

- **Continuous Bag-of-Words** (CBOW),
- **Skip-Gram** (SG).

These approaches are historically interesting, but now a bit outdated... ;-)

<sup>1</sup> Tomas Mikolov & all, (2013), [W3VEC]

CBOW : Continuous Bag of Words - Embedding based on the prediction of the word according to its context.

SG : Skip-gram - Embedding based on context prediction from the word.



(french)

## Word2Vec<sup>1</sup>

Approche ayant pour objectif de constituer des **dictionnaires** dont la représentation vectorielle des mots est basée sur le **contexte** et donc de la **sémantique**.

Des dictionnaires construits à partir de gros corpus sont disponibles.

Deux modèles :

- **Continuous Bag-of-Words** (CBOW),
- **Skip-Gram** (SG).

Ces approches sont historiquement intéressantes, mais désormais un peu dépassées... ;-)

<sup>1</sup> Tomas Mikolov & all, (2013), [W3VEC]

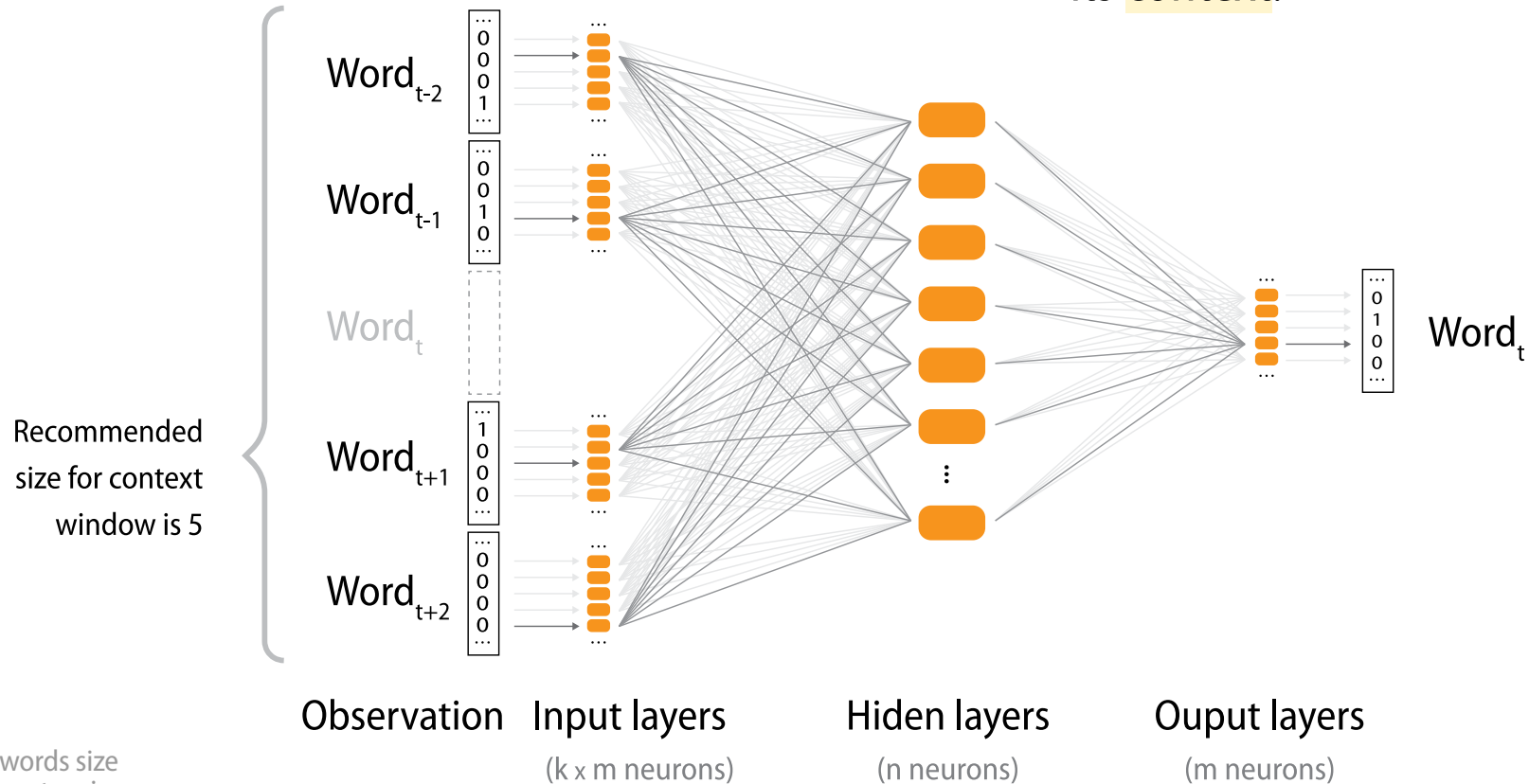
CBOW : Continuous Bag of Words - Embedding based on the prediction of the word according to its context.

SG : Skip-gram - Embedding based on context prediction from the word.



# Continuous Bag-of-Words

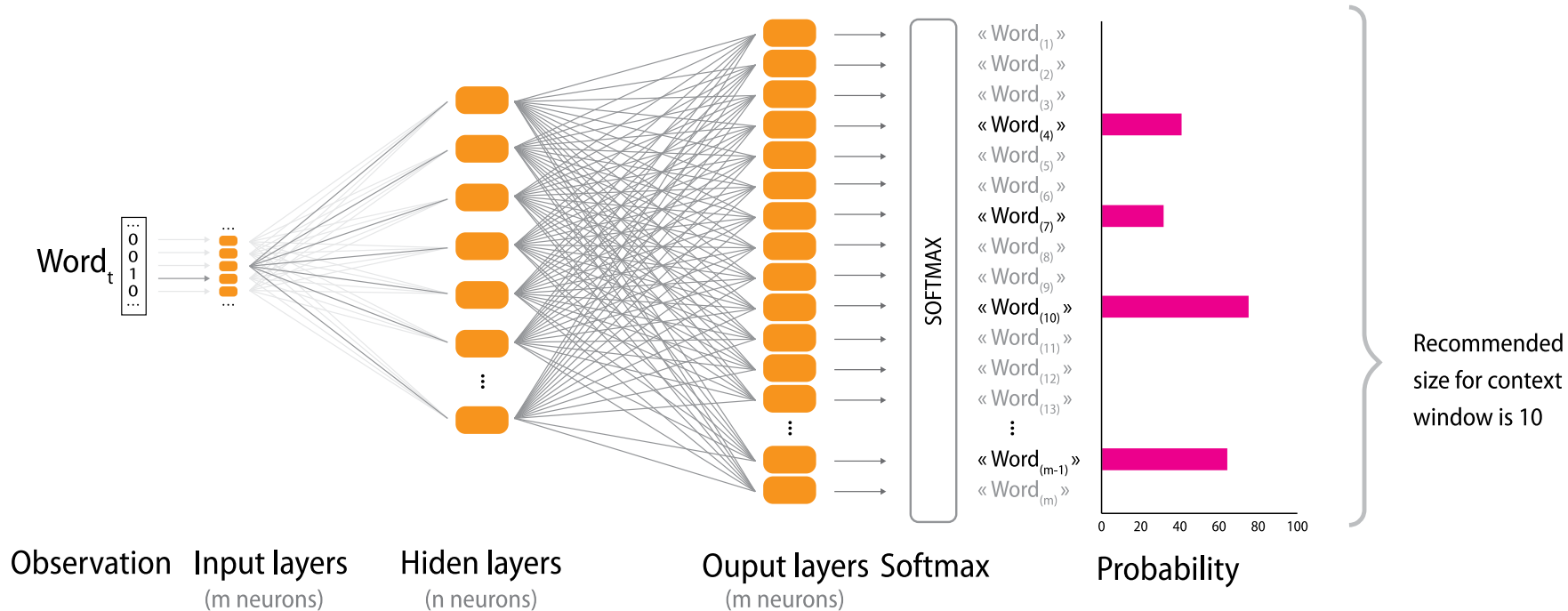
” The objective of training is to (re)find a word from its context.



$k$  = Bag of words size  
 $n$  = Dense vector size  
 $m$  = Dictionary size



The objective of training is to (re)find the context from the word.



k = Bag of words size  
n = Dense vector size  
m = Dictionary size

## GloVe<sup>1</sup>



Contrairement à Word2vec, GloVe ne repose pas uniquement sur des statistiques locales (informations sur le contexte local des mots), mais intègre des **statistiques globales** (co-occurrence des mots) pour obtenir des vecteurs de mots.

<sup>1</sup> Jeffrey Pennington & all, (2014), [GLOVE]  
Training is performed on aggregated global word-word co-occurrence statistics

## (Flau)BERT<sup>1</sup>

2018



BERT est une représentation du langage proposée par Google en 2018, permettant de prendre en compte la dimension contextuelle du langage (« Avocat », peut être un fruit ou un juriste..).

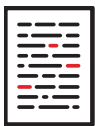
FlauBERT<sup>2</sup> est une adaptation de l'algorithme au français.

To be continued !

<sup>1</sup> BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.  
Jacob Devlin, Ming-Wei Chang, Kenton Lee,  
Kristina Toutanova  
<https://arxiv.org/abs/1810.04805>

<sup>2</sup> FlauBERT: Unsupervised Language Model Pre-training for French.  
Hang Le, Loïc Vial, Jibril Frej, Vincent Segonne, Maximin Coavoux,  
Benjamin Lecouteux, Alexandre Allauzen, Benoît Crabbé,  
Laurent Besacier, Didier Schwab  
<https://arxiv.org/abs/1912.05372>

5b



**Sparse data (text)**  
Embedding

5.1

## Text encoding

- From text to tensor
- One hot encoding
- Embedding

5.2

## Example 1 : IMDB1

- Sentiment analysis with one-hot encoding



5.3

## Example 2 : IMDB2/3/4

- Sentiment analysis with embedding

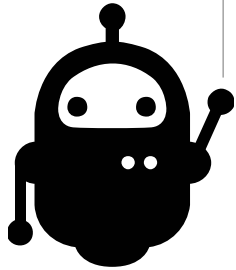






# One-Hot encoding with IMDB

Notebook : [\[IMDB1\]](#)



## **Objective :**

Guess whether a film review is positive or not based on the analysis of the text, using One-Hot encoding.

## **Dataset :**

The IMDB dataset is composed of 50,000 film reviews from the site of the same name.

<https://www.imdb.com/>





<60 s

Import  
and init

START

Step 1

Understanding  
hot-one  
encoding

Step 2

Step 3

About our  
dataset

Step 4

Step 5

Build the  
model

Step 6

Step 7

Evaluate

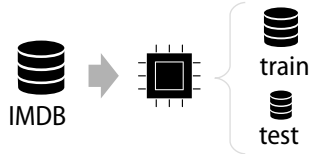
Step 8

END

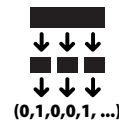
## Objectives :

Analysis of  
film reviews  
with  
One-Hot  
encoding

Retrieve data



One-hot  
vector  
encoding



Train the  
model



5b



**Sparse data (text)**  
Embedding

5.1

## Text encoding

- From text to tensor
- One hot encoding
- Embedding

5.2

## Example 1 : IMDB1

- Sentiment analysis with one-hot encoding



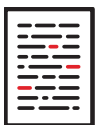
5.3

## Example 2 : IMDB2/3/4

- Sentiment analysis with embedding



5b



**Sparse data (text)**  
Embedding

5.1

## Text encoding

- From text to tensor
- One hot encoding
- Embedding

5.2

## Example 1 : IMDB1

- Sentiment analysis with one-hot encoding



5.3

## Example 2 : IMDB2/3/4

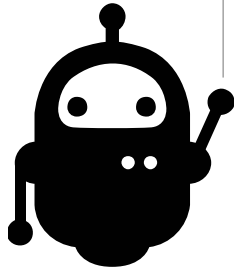
- Sentiment analysis with embedding





# Text embedding with IMDB

Notebook : [IMDB2-4]



## **Objective :**

Guess whether a film review is positive or not based on the analysis of the text, using embedding.

## **Dataset :**

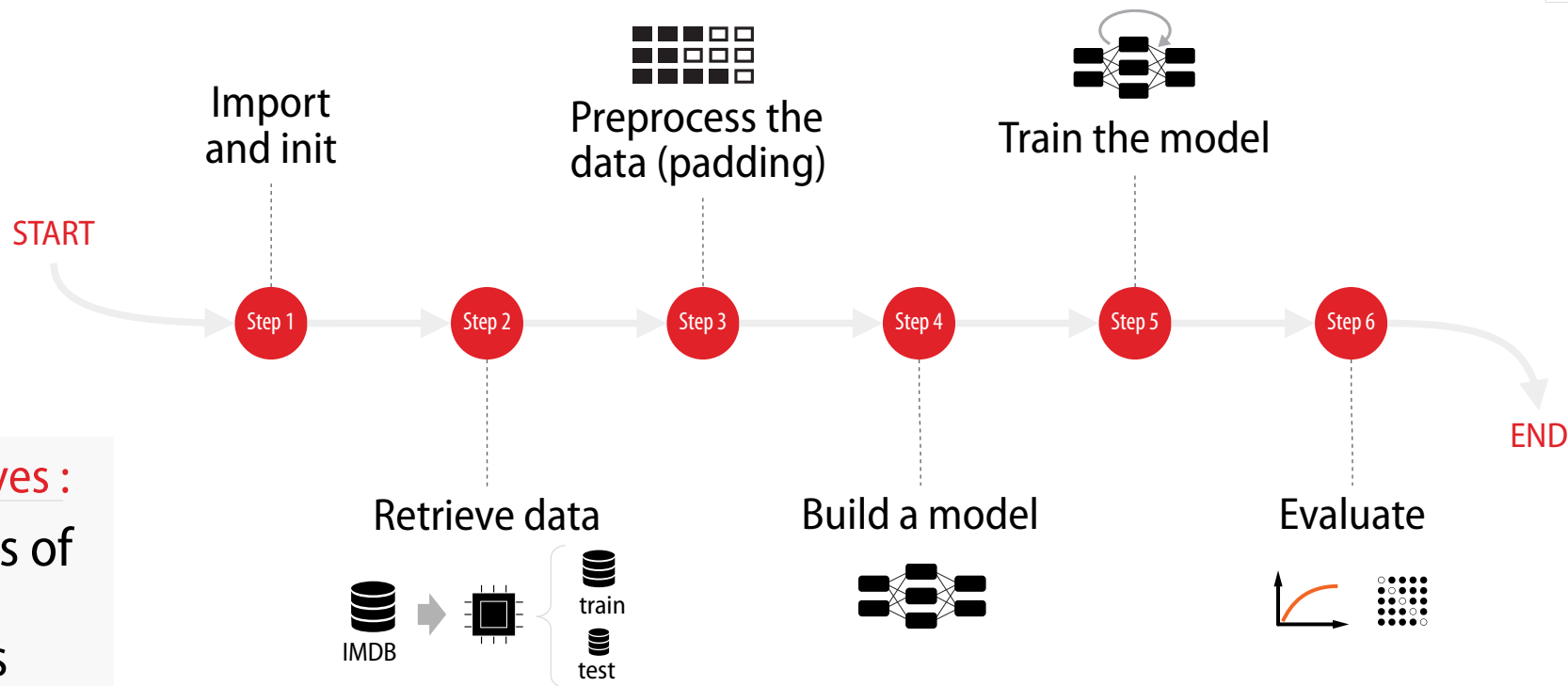
The IMDB dataset is composed of 50,000 film reviews from the site of the same name.

<https://www.imdb.com/>





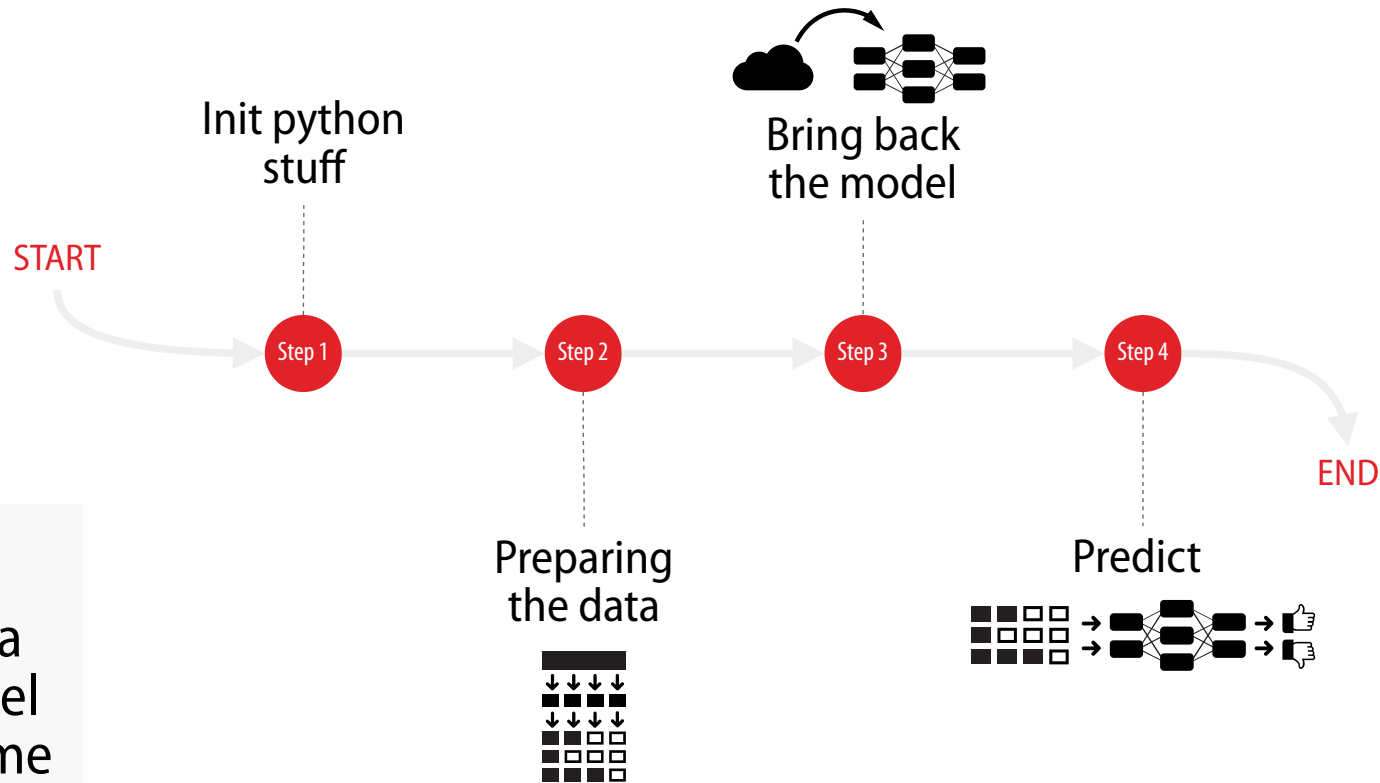
<60 s



## Objectives :

Analysis of  
film  
reviews  
with Keras  
embedding

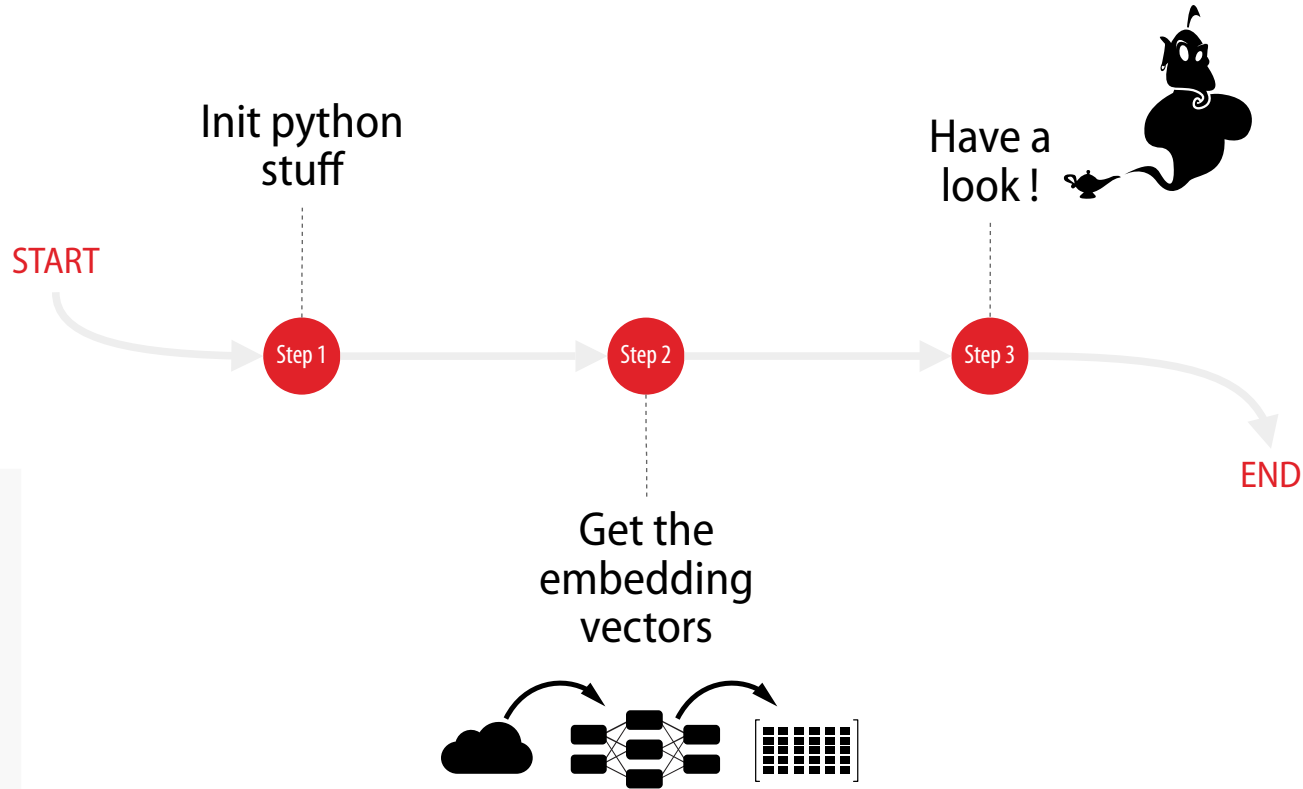




## Objectives :

Retrieving a saved model to treat some new reviews





## Objectives :

Retrieving a saved model and play with embedded vectors





# Next, on Fidle :

6



**Sequences data**  
RNN



## Jeudi 5 janvier,

Épisode 6 :

### **Quand les données sont des séquences...**

Données séquentielles et réseaux récurrents (RNN)  
RNN - LSTM - GRU - Spécificités des données séquentielles.

Exemples proposés :

Tentative de prédiction de la trajectoire d'une coccinelle (virtuelle)  
et de la météorologie, à partir de données réelles, à 3h et 12h.

Durée : 2h00



Next on Fidle :

6



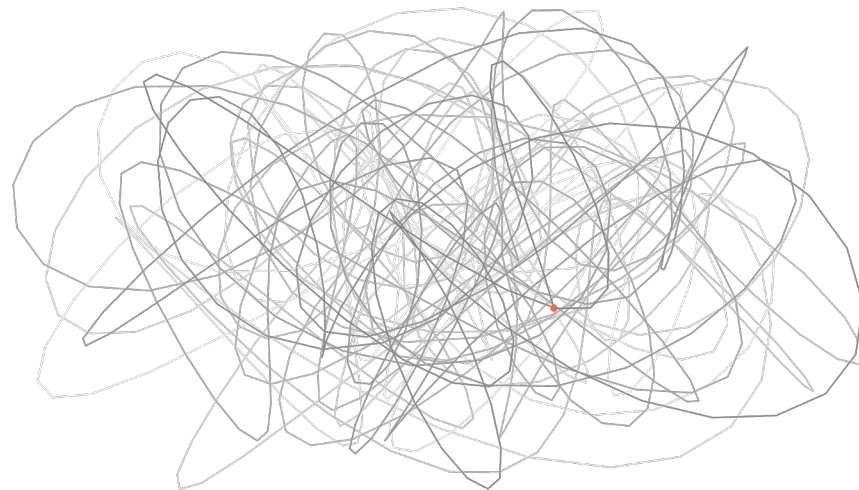
Sequence data  
RNN



**Jeudi 20 janvier,**

Séquence 6 :

**Quand les données sont des séquences...  
Données séquentielles et réseaux  
récurrents (RNN)**



To be continued...



Contact@fidle.cnrs.fr



FIDLE <https://fidle.cnrs.fr>



YouTube <https://fidle.cnrs.fr/youtube>



Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0)  
<https://creativecommons.org/licenses/by-nc-nd/4.0/>