

# Deep Learning pour le Text Mining

## L'algorithme « word2vec »

Ricco Rakotomalala

1. Prise en compte du contexte
2. Les algorithmes SKIP-GRAM et CBOW
3. Représentation des documents
4. Bibliographie

Insuffisances de la représentation usuelle en « sac de mots »

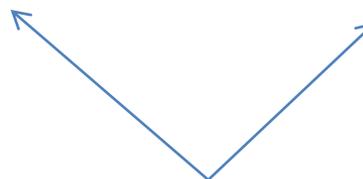
**PRISE EN COMPTE DU CONTEXTE EN TEXT MINING**

La représentation en sac de mots ne tient pas compte des positions relatives de mots dans les documents.

- (0) condition du bien etre
- (1) etre important
- (2) solution bien etre
- (3) important bien etre



	bien	condition	du	etre	important	solution
(0)	1	1	1	1	0	0
(1)	0	0	0	1	1	0
(2)	1	0	0	1	0	1
(3)	1	0	0	1	1	0



La contextualisation de « bien » par « être » (ou l'inverse d'ailleurs) est importante, ils sont souvent « voisins ». La représentation en sac de mots passe à côté (les algos de machine learning ne verront que la co-occurrence, c'est déjà pas mal - ex. topic modeling)

Idée du prolongement lexical : déterminer une représentation des termes par un vecteur numérique de dimension  $K$  (paramétrable), en tenant compte de son contexte (fenêtre de voisinage  $V$  dont la taille est paramétrable).

(0) condition du bien etre  
(1) etre important  
(2) solution bien etre  
(3) important bien etre

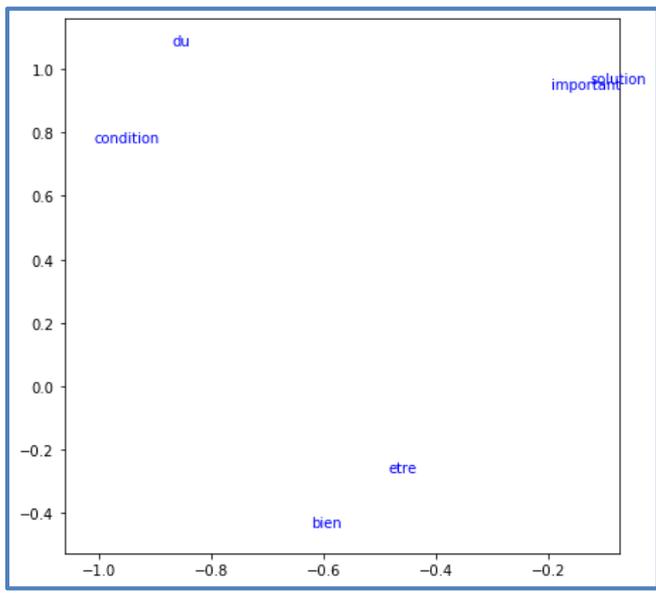


Les termes apparaissant dans des contextes similaires sont proches (au sens de la distance entre les vecteurs de description).



A partir de la représentation des termes qui les composent, il est possible de dériver une description numérique (vectorielle) des documents.

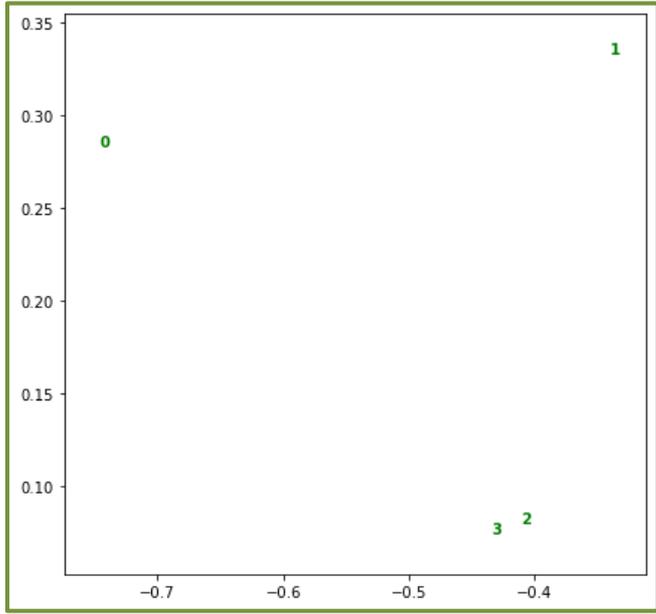
- (0) condition du bien etre
- (1) etre important
- (2) solution bien etre
- (3) important bien etre



### Représentation des termes

*Bon, converger sur 4 observations n'est pas évident quoiqu'il en soit*

- (A)  $K \ll$  taille (dictionnaire), réduction forte de la dimensionnalité
- (B) Il est possible d'effectuer des traitements de machine learning à partir de ce nouvel espace de représentation (clustering, classement,...)



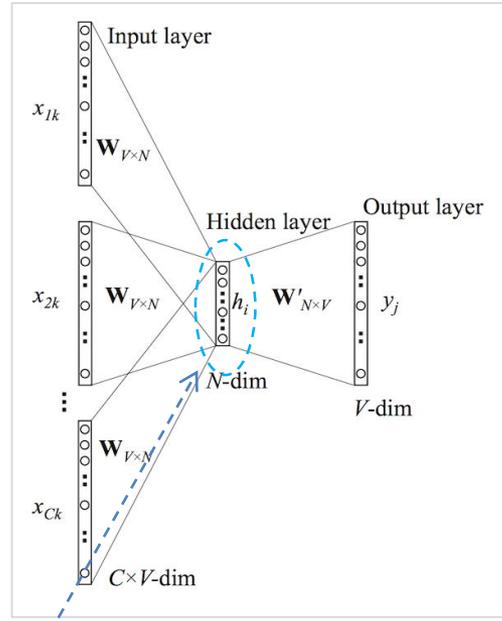
### Représentation des documents

Deep Learning pour le prolongement lexical (word embedding)

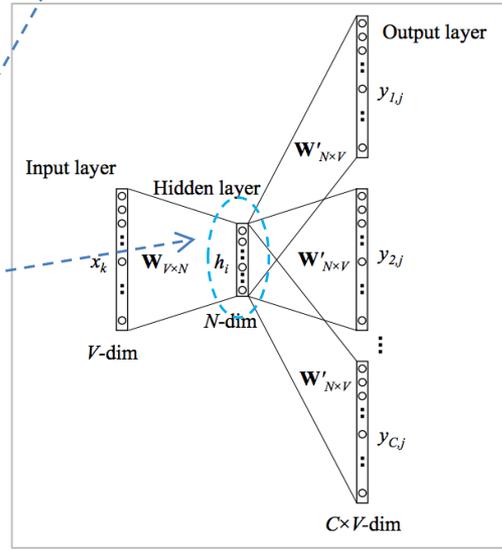
# L'ALGORITHME WORD2VEC

Modéliser les termes en utilisant un réseau de neurones (un perceptron à une couche cachée) avec en entrée le contexte (le voisinage) et en sortie le terme (CBOW) ou inversement (SKIP-GRAM).

A la manière des auto-encodeurs, ce sont les descriptions à la sortie de la couche cachée qui nous intéressent (nouvelles coordonnées des termes). Elles constituent la **représentation des termes dans un nouvel espace**.



CBOW

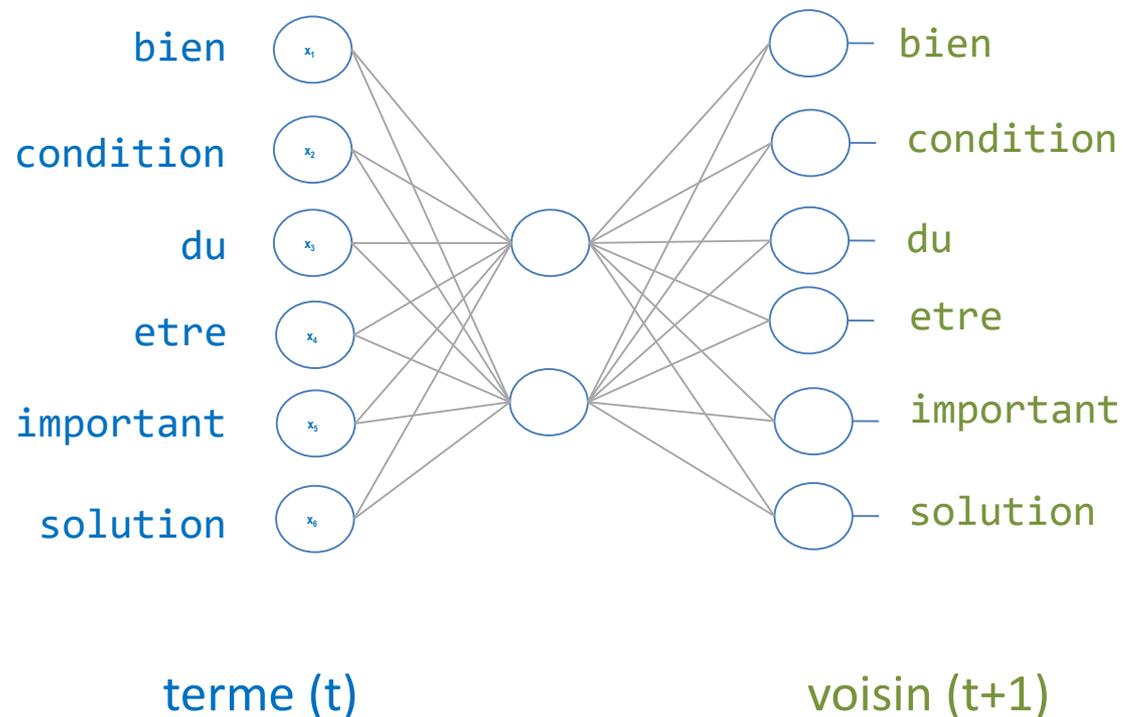


SKIP-GRAM

[https://fr.wikipedia.org/wiki/Word\\_embedding](https://fr.wikipedia.org/wiki/Word_embedding)

Modéliser les voisinages à partir des termes c.-à-d.  $P(\text{voisin}[s] / \text{terme})$ .

Ex. le voisin immédiat qui succèdent les termes dans les documents



L'astuce passe par un encodage approprié des données tenant compte du voisinage. Ex. voisinage de taille 1 vers l'avant  $v_{t+1}$

Description BOW (bag of words)

	bien	condition	du	etre	important	solution
condition du bien etre	1	1	1	1	0	0
etre important	0	0	0	1	1	0
solution du bien etre	1	0	0	1	0	1
important bien etre	1	0	0	1	1	0

**Entrée (terme t)**

Terme	bien	condition	du	etre	important	solution
condition	0	1	0	0	0	0
du	0	0	1	0	0	0
du	0	0	1	0	0	0
bien	1	0	0	0	0	0
bien	1	0	0	0	0	0
bien	1	0	0	0	0	0
etre	0	0	0	1	0	0
important	0	0	0	0	1	0

*etc.*



**Sortie (voisin t + 1)**

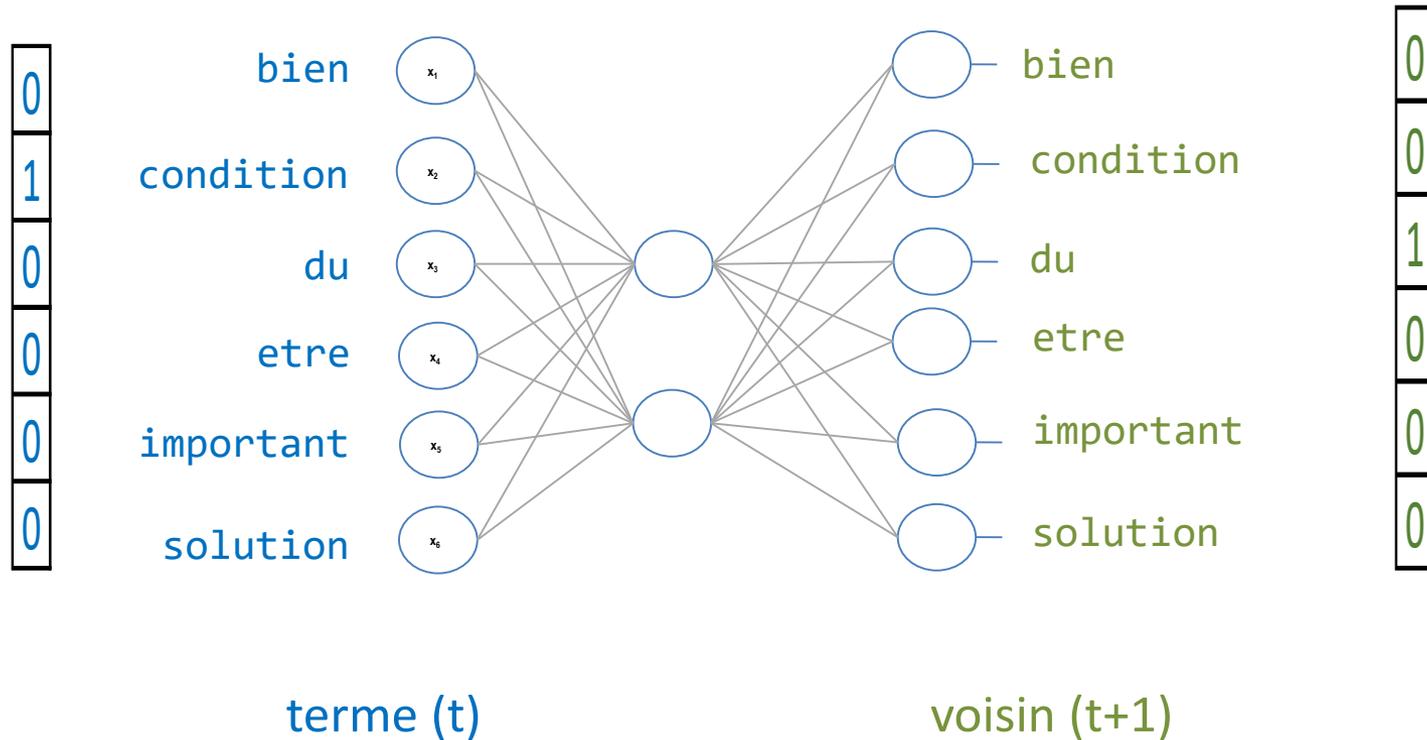
Terme	bien	condition	du	etre	important	solution
du	0	0	1	0	0	0
bien	1	0	0	0	0	0
bien	1	0	0	0	0	0
etre	0	0	0	1	0	0
etre	0	0	0	1	0	0
etre	0	0	0	1	0	0
important	0	0	0	0	1	0
bien	1	0	0	0	0	0

*etc.*

Ce sont ces données (pondération forcément binaire) que l'on présentera au réseau.  
On est dans un (une sorte de) schéma d'apprentissage supervisé multi-cibles



Exemple d'une observation présentée au réseau : (condition → du)



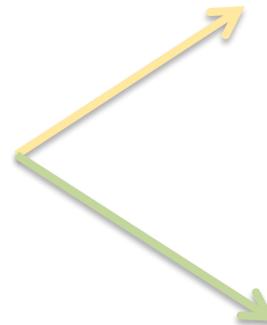
# Double tableau pour la sortie maintenant : voisinages (t-1) et (t+1)

	bien	condition	du	etre	important	solution
condition du bien etre	1	1	1	1	0	0
etre important	0	0	0	1	1	0
solution du bien etre	1	0	0	1	0	1
important bien etre	1	0	0	1	1	0

**Entrée (terme t)**

Terme	bien	condition	du	etre	important	solution
du	0	0	1	0	0	0
du	0	0	1	0	0	0
bien	1	0	0	0	0	0
bien	1	0	0	0	0	0
bien	1	0	0	0	0	0

etc.



**Sortie (voisin t -1)**

Terme	bien	condition	du	etre	important	solution
condition	0	1	0	0	0	0
solution	0	0	0	0	0	1
du	0	0	1	0	0	0
du	0	0	1	0	0	0
important	0	0	0	0	1	0

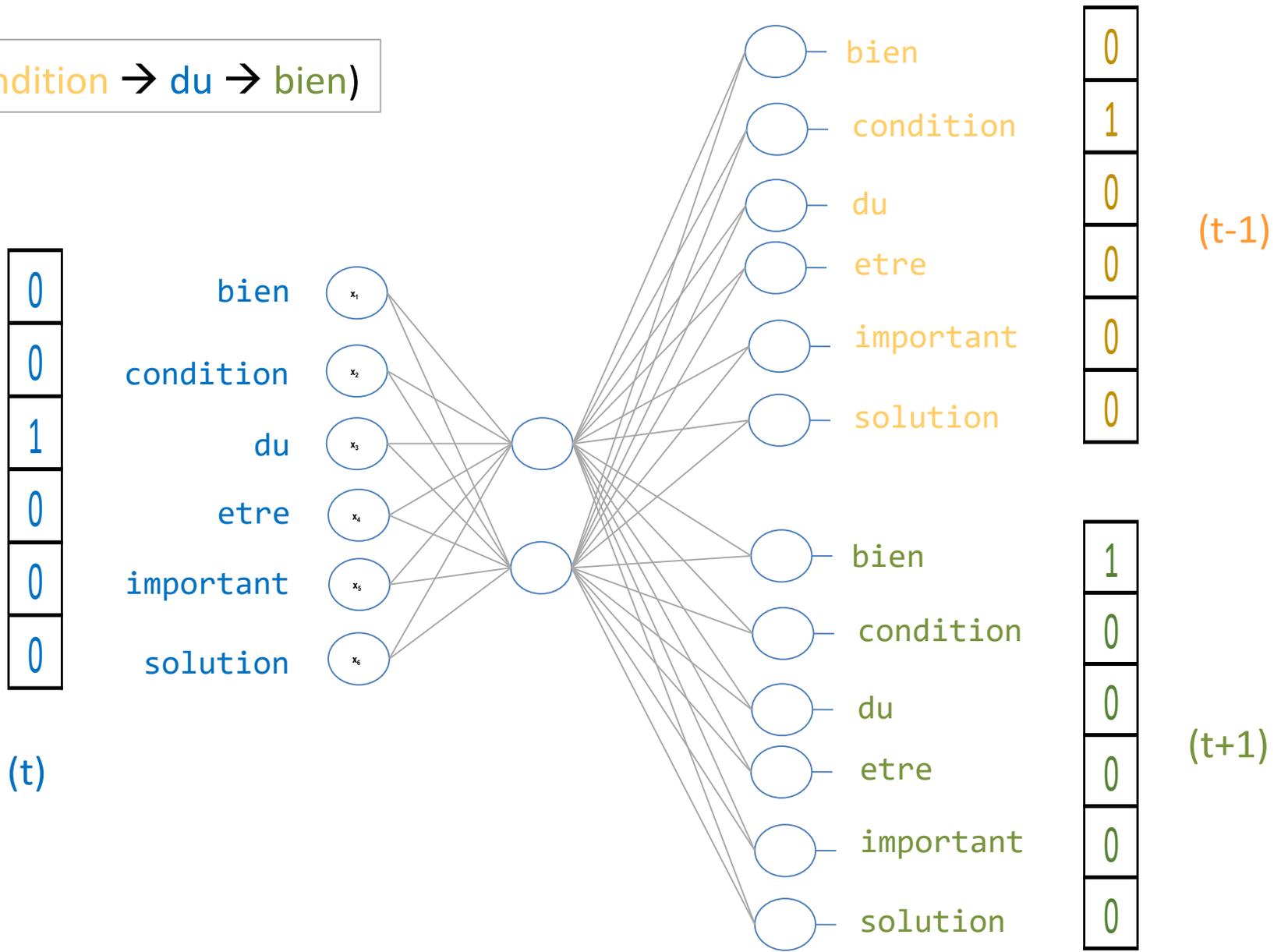
etc.

**Sortie (voisin t + 1)**

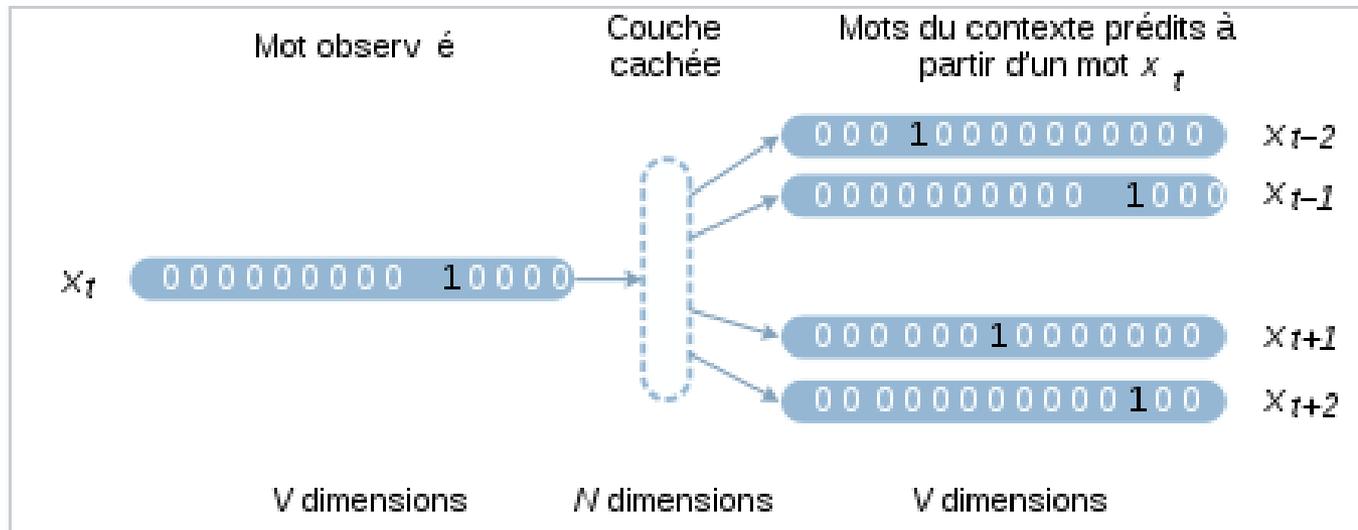
Terme	bien	condition	du	etre	important	solution
bien	1	0	0	0	0	0
bien	1	0	0	0	0	0
etre	0	0	0	1	0	0
etre	0	0	0	1	0	0
etre	0	0	0	1	0	0

etc.

Ex. (condition → du → bien)

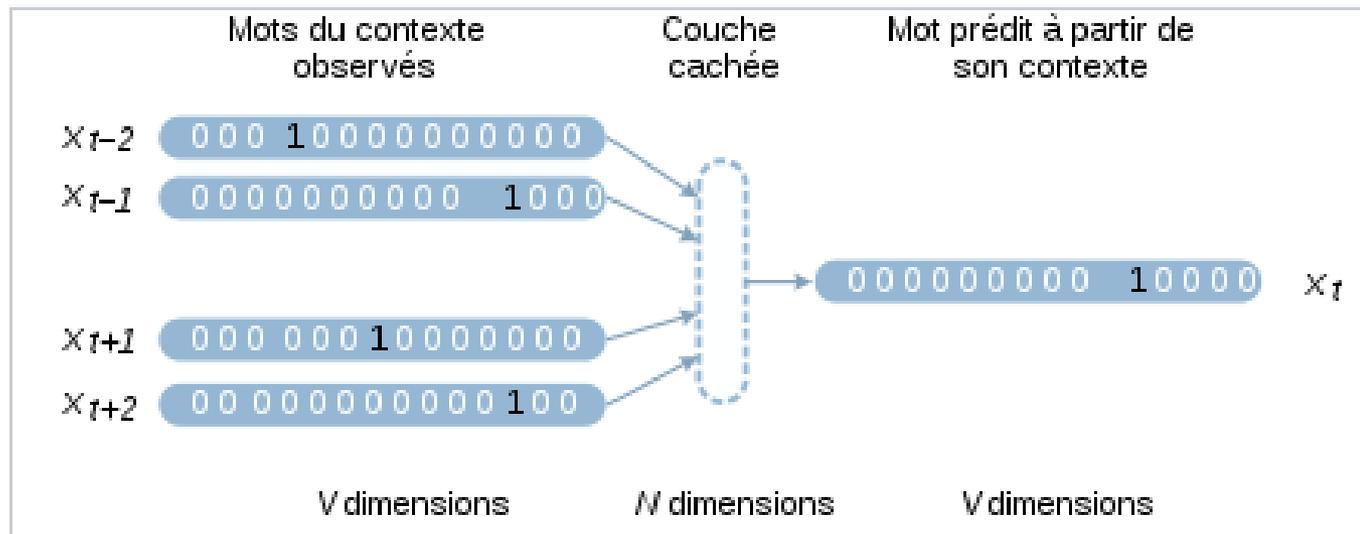


Il est possible de prendre un voisinage plus étendu ( $V = 2$  ou plus). Attention simplement à la dilution de l'information.



[https://fr.wikipedia.org/wiki/Word\\_embedding](https://fr.wikipedia.org/wiki/Word_embedding)

La problématique est inversée : on s'appuie sur le voisinage (le contexte) pour apprendre les termes. On modélise  $P(\text{terme} / \text{voisin}[s])$ .



[https://fr.wikipedia.org/wiki/Word\\_embedding](https://fr.wikipedia.org/wiki/Word_embedding)

- La fonction de transfert pour la couche centrale est linéaire
- Pour la couche de sortie, la fonction de transfert est [softmax](#)
- La « [negative log-likelihood](#) » fait office de fonction de perte (à la place du classique MSE – mean squared error). En conjonction avec softmax, le calcul du gradient en est [largement simplifiée](#) lors de l'optimisation
- Dixit la [documentation de H2O](#) (fameux package de Deep Learning) : skip-gram donne plus de poids aux voisins proches, elle produit de meilleurs résultats pour les termes peu fréquents

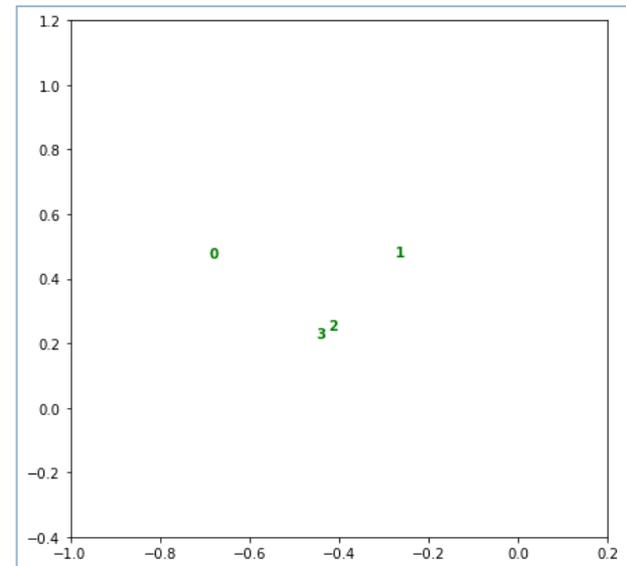
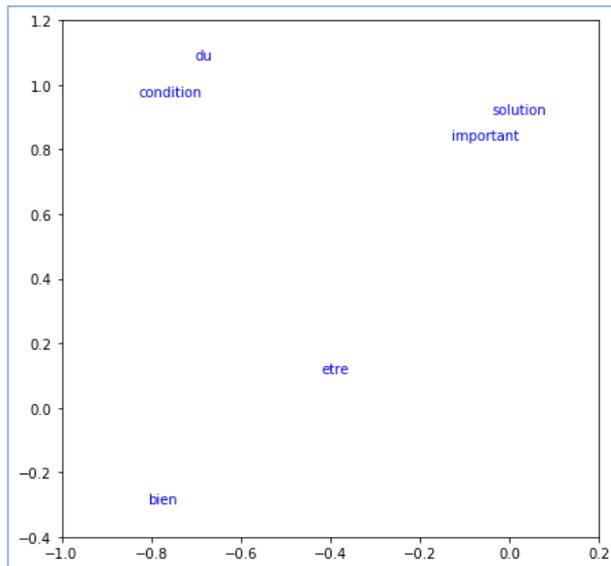
Passer des termes aux documents : vecteur de représentation

# REPRÉSENTATION DES DOCUMENTS

Disposer d'une représentation des documents dans le nouvel espace est indispensable pour pouvoir appliquer les algorithmes de machine learning (ex. catégorisation de documents, etc.)

```
<document>  
< sujet>acq</sujet>  
<texte>  
Resdel Industries Inc said it has agreed to acquire San/Bar Corp in a share-for-share exchange, after San/Bar distributes all shgares of its Break-Free Corp subsidiary.  
</texte>  
</document>  
<document>  
< sujet>acq</sujet>  
<texte>  
Warburg, Pincus Capital Co L.P., an investment partnership, said it told representatives of Symbion Inc it would not increase the 3.50-dlr-per-share cash price it has offered for the company.  
</texte>  
</document>
```

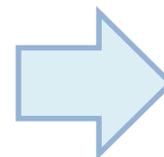
Comment passer de la représentation des termes à celle des documents (composés de termes) ?



Solution simple : calculer la moyenne des coordonnées (barycentre) des termes qui composent le document.

- (0) condition du bien etre
- (1) etre important
- (2) solution bien etre
- (3) important bien etre

	Word	V1	V2
0	du	-0.703333	1.077643
1	condition	-0.828805	0.966393
2	solution	-0.040079	0.911001
3	important	-0.127894	0.830322
4	bien	-0.808076	-0.295098
5	etre	-0.419983	0.109700



	C1	C2
(0)	-0.690049	0.464660
(1)	-0.273938	0.470011
(2)	-0.422713	0.241868
(3)	-0.451984	0.214975

Calcul des coordonnées du document n°1

$$C1(1) = \frac{-0.419983 + (-0.127894)}{2} = -0.273938$$

$$C2(1) = \frac{0.109700 + 0.830322}{2} = 0.470011$$

Si (K = 2), une représentation simultanée dans le plan est possible.



# CONCLUSION

- « word2vec » est une technique de « word embedding » basée sur un algorithme de deep learning.
- L'objectif est de représenter les termes d'un corpus à l'aide d'un vecteur de taille  $K$  (paramètre à définir, parfois des centaines, tout dépend de la quantité des documents), où ceux qui apparaissent dans des contextes similaires (taille du voisinage  $V$ , paramètre à définir) sont proches (au sens de la dist. cosinus par ex.).
- De la description des termes, nous pouvons dériver une description des documents, toujours dans un espace de dimension  $K$ . Possibilité d'appliquer des méthodes de machine learning par la suite (ex. catégorisation de documents).
- $K \ll$  taille du dictionnaire : nous sommes bien dans la réduction de la dimensionnalité (par rapport à la représentation « bag-of-words » par ex.).
- Il existe des modèles pré-entraînés sur des documents (qui font référence, ex. Wikipedia ; en très grande quantité) que l'on peut directement appliquer sur nos données (ex. [Google Word2Vec](#) ; [Wikipedia2Vec](#))

Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. “[Efficient Estimation of Word Representations in Vector Space.](#)” In Proceedings of Workshop at ICLR. (Sep 2013)

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. “[Distributed Representations of Words and Phrases and their Compositionality.](#)” In Proceedings of NIPS. (Oct 2013)

H2O, “Word2Vec”, <http://docs.h2o.ai/h2o/latest-stable/h2o-docs/data-science/word2vec.html>