

Application :

Analyse de sentiment  
Traduction  
Reconnaissance vocale

Pourquoi adapter au français ?

Jamais fait actuellement Algo anglais appliqué

Pas adaptée en :  $\mathbb{C}$ ,  $\mathbb{R}$ ,  $\mathbb{Z}$  ou  $\mathbb{N}$  <sup>personne</sup>

~~Nuit~~ au } Nos sommes <sup>4</sup>is  
Traduction } \* EX plus simple d'initiation.

~~De nos jours~~, la compréhension de notre langage (c-à-d le français)  
par les ordinateurs est de plus en plus important ; car ils sont de plus en plus  
présent dans notre vie ~~et le seront de plus en plus~~ et le seront de + en +, ex :  
Smartphon, Voiture Autonome, ~~Robot Siri~~ "Siri".

Quelques .



## Le français dans l'ère du "Deep Learning" et les "Word Embeddings"

Ce projet consiste en l'utilisation de techniques connues avec le nom de "Word Embeddings" pour le français. Ces techniques ont été largement étudiées pour l'anglais [1][2][3][4] et récemment pour des autres langues. Normalement, les auteurs de cette technique ont démontré qu'elle avait une meilleure performance par des évaluations des analogies [3][4]. C'est-à-dire, qu'une méthode était considérée comme supérieure si elle arrivait à trouver correctement un des mots d'une analogie avec 4 mots. Par exemple, la tâche consiste à trouver correctement *Queen* sachant qu'on connaît *Man*, *King* et *Woman*. Presque 20 mill. exemples sont données par les auteurs de cette technique [3]. Cependant, l'évaluation est appliquée que pour l'anglais et complètement inexistant pour le français ou autres langues. Également, la source avec laquelle la méthode est entraînée est sans aucune doute un ressource précieux pour trouver de bons résultats. Pour l'anglais, la Wikipédia a été utilisé avec de bons résultats. Pour le français, on se contente d'utiliser la collection objective et pas une collection générique afin de comprendre les caractéristiques particulières de la langue. Probablement la manque d'une ressource d'évaluation est à l'origine de cette situation.

Les principaux objectifs de ce projet sont :

- La traduction et analyses des collections d'évaluation par analogie existants pour l'évaluation de words embeddings pour la langue française.
- La construction de modèles avec les méthodes existants et de ressources ouverts pour la langue française.
- L'évaluation et comparaison de différentes implémentations et configurations avec la collection d'évaluation par analogie en français.

Plusieurs implémentations sont disponibles :

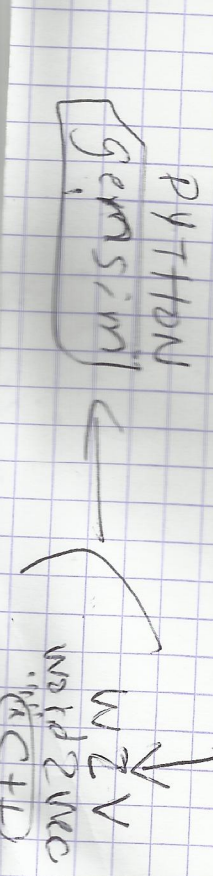
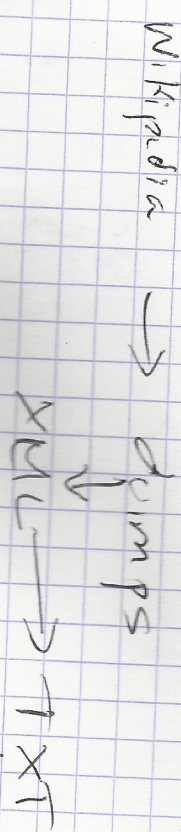
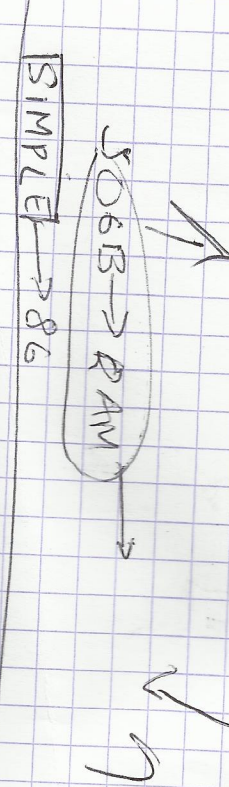
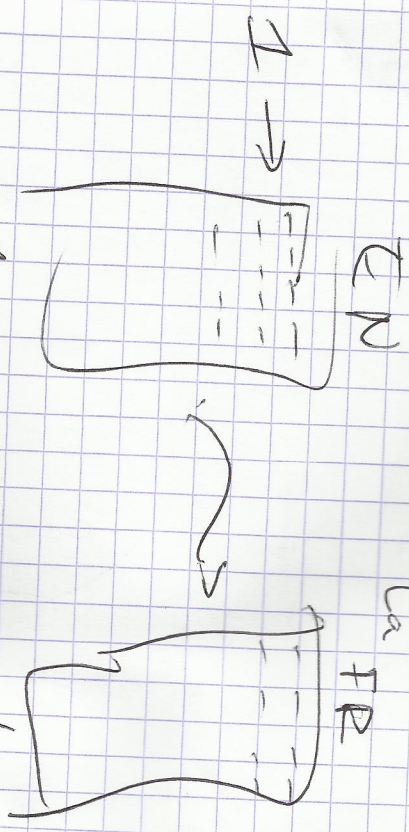
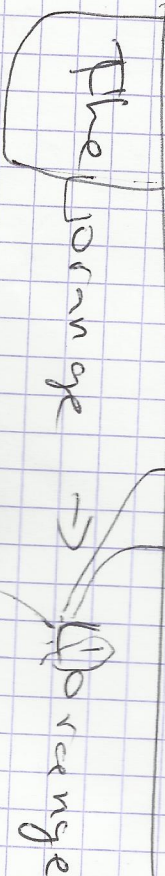
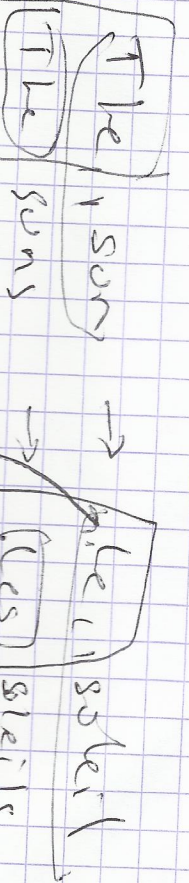
- Gensim <https://radimrehurek.com/gensim/>
- hyperwords <https://bitbucket.org/omerlevy/hyperwords>
- - word2vec <https://code.google.com/archive/p/word2vec/>
- tensorflow <https://www.tensorflow.org/tutorials/word2vec/>

Procédure recommandée :

- Télécharger les fichiers questions-words.txt et question-phrases.txt disponibles sur <https://code.google.com/archive/p/word2vec/> (chercher autres collections pertinents)
- Trouver les mots uniques dans chaque fichier, les traduire et faire le mapping
- Vérifier que les analogies sont toujours valides (pas de problèmes de traduction) et proposer des analogies équivalents pour le français
- Télécharger la Wikipédia en français (fichier frwiki-20170101-pages-articles.xml.bz2 sur <https://dumps.wikimedia.org/frwiki/20170101/> ou le plus récent)
- Traiter la Wikipédia avec les implémentations proposées
- Utiliser les vecteurs obtenus pour calculer la performance avec chaque combinaison des paramètres (au moins deux calculs peuvent être faites, 3COSADD et 3COSMUL [1])
- Rédiger un rapport avec les résultats et les conclusions de ce projet

- [1] Omer Levy and Yoav Goldberg. Linguistic Regularities in Sparse and Explicit Word Representations.
- [2] Omer Levy, Yoav Goldberg, Ido Dagan. Improving distributional similarity with lessons learned from word embeddings.
- [3] T Mikolov, I Sutskever, K Chen, GS Corrado, J Dean. Distributed representations of words and phrases and their compositionality
- [4] T Mikolov, K Chen, G Corrado, J Dean. Efficient estimation of word representations in vector space



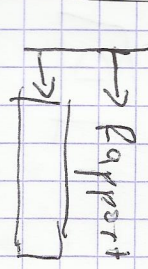


French W2V

IRF M21

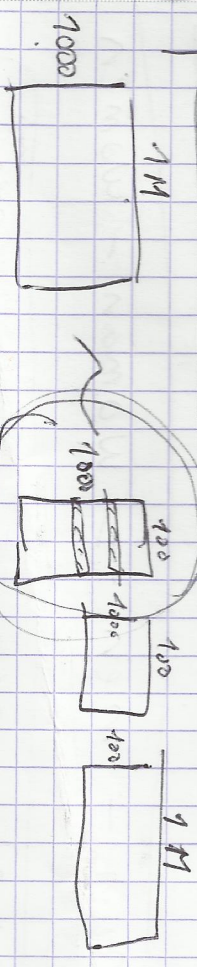
6/4 → 1/4/2017

04/01/2017



2013 → Word Embeddings

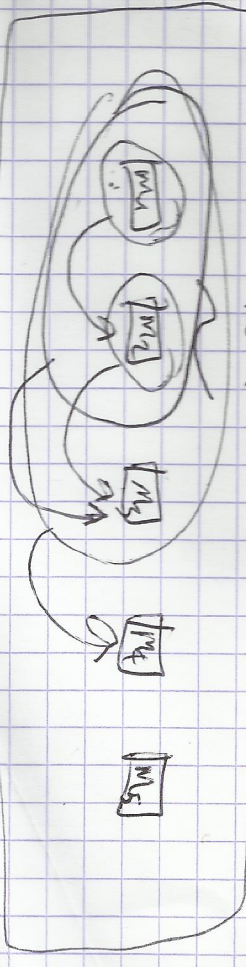
	$d_1$	$d_2$	...	$d_n$
$m_1$	1	1		
$m_2$	0	1		1
$m_3$	3	0		
$m_i$	0	1		1



2013 M1 Kolar

RN

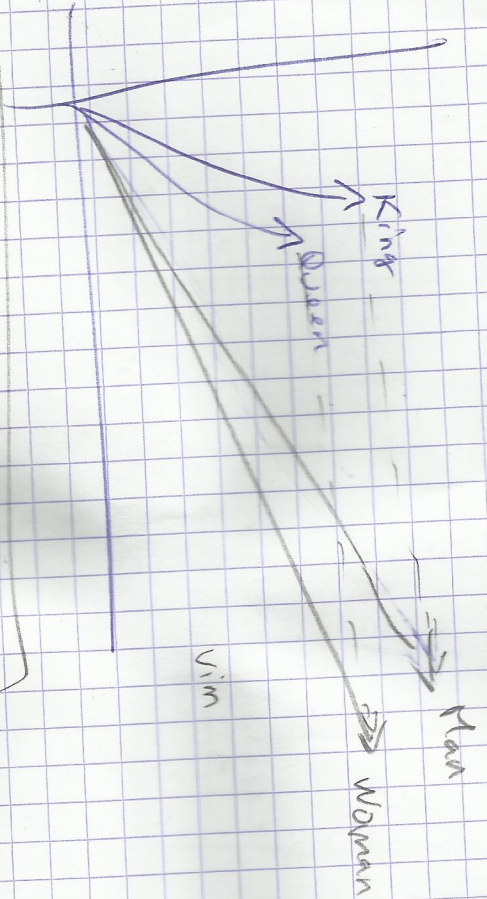
LSI → SVD





King - Men & Queen - Woman

Pages  
Titles



King - Man + Woman ~ Queen

Carbure  
The tree

don't

all +

King - Man = Queen - Woman

A + B - C = D

A + C = D - B

King - Man = Queen - Woman

King - Man + Woman ~ Queen

Man - men = Woman - women

Man - men = pen - pens



3h-6h

10 11 12  
13 14 15  
16 17 18

Wikipedia

	W = 5		W = 50	
	Stigman	CBOW	Stigman	CBOW
English	40%			
French	30%	20%	25%	58%

Link

Link