

# Méthodes d'évaluation pour la substitution de vecteurs de mots

Stanislas Morbieu <sup>1,2</sup>, François Role <sup>1</sup>, Mohamed Nadif <sup>1</sup>

<sup>1</sup> LIPADE, Université Paris Descartes

<sup>2</sup> Kernix

XXVI èmes Rencontres de la Société Francophone de  
Classification

Représentation vectorielle de mots

Mots hors vocabulaire

Similarité entre les vecteurs originaux et les substituts

Comportement sur des tâches de similarité et d'analogie

## Idée principale

"You shall know a word by the company it keeps." John R. Firth (1957)

## Idée principale

"You shall know a word by the company it keeps." John R. Firth (1957)

Je visite la place\_Stanislas à Nancy .

mot cible

contexte

Je visite la place\_Stanislas à Nancy .

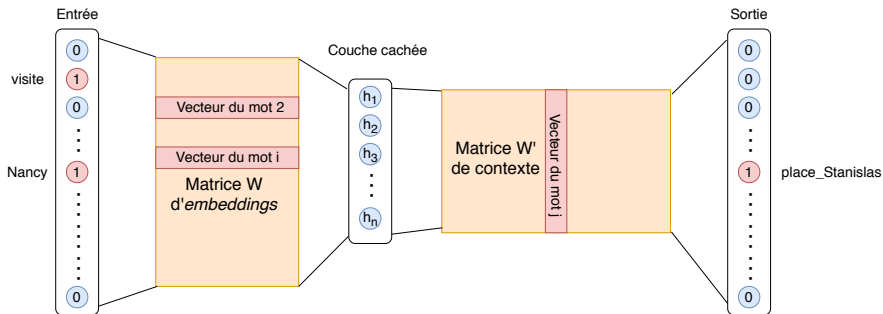


Figure 1: Word2vec CBOW [Mikolov 2013]

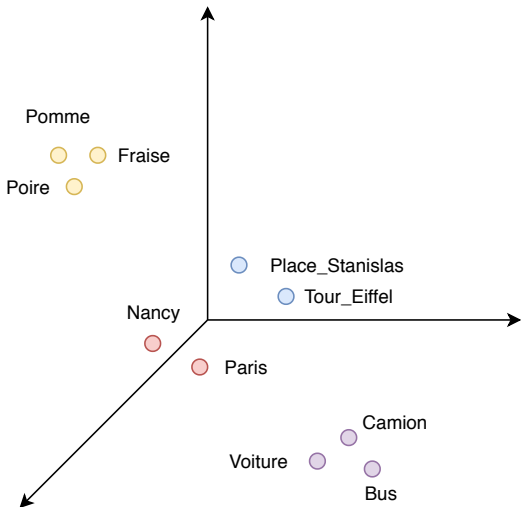


Figure 2: Représentation des mots dans l'espace vectoriel.

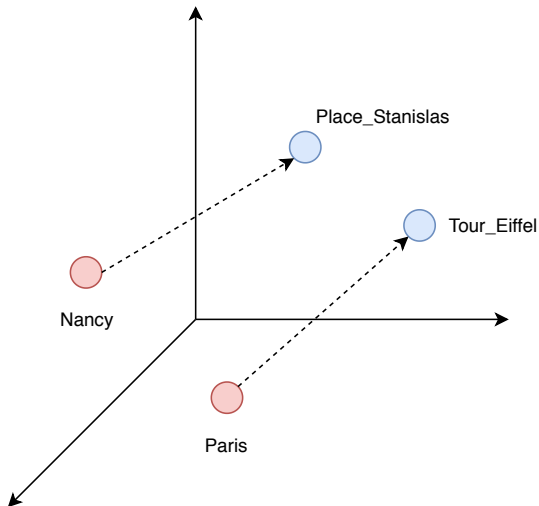
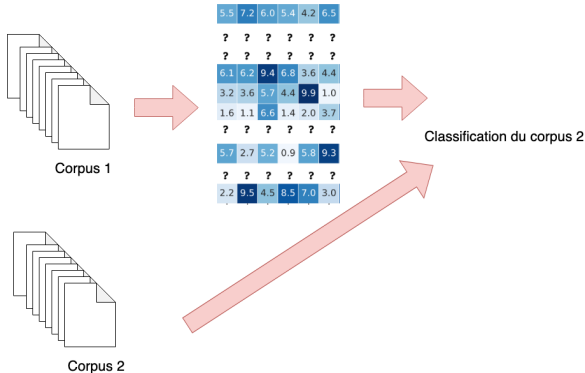


Figure 3: Représentation des analogies dans l'espace vectoriel.



## Ordre de grandeur

Près de 61% des mots du jeu de données NG20<sup>a</sup> sont absents dans le modèle Word2vec entraîné sur le corpus Google News.

<sup>a</sup><http://qwone.com/~jason/20Newsgroups/>



# Gestion des mots hors vocabulaire

## Différentes méthodes :

- Initialiser les vecteurs inconnus avec des **zéros** ou des valeurs **aléatoires**.

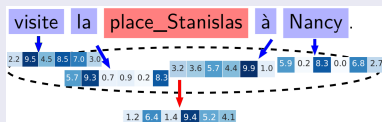
$$[?, ?, ?, \dots, ?] \rightarrow [0, 0, 0, \dots, 0]$$

# Gestion des mots hors vocabulaire

## Différentes méthodes :

- Initialiser les vecteurs inconnus avec des **zéros** ou des valeurs **aléatoires**.
- Prendre le **centroïde des vecteurs des mots** d'une fenêtre de contexte.

$$[?, ?, ?, \dots, ?] \rightarrow [0, 0, 0, \dots, 0]$$

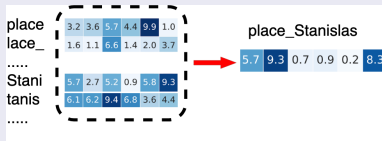
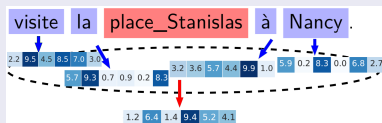


# Gestion des mots hors vocabulaire

## Différentes méthodes :

- Initialiser les vecteurs inconnus avec des **zéros** ou des valeurs **aléatoires**.
- Prendre le **centroïde** des **vecteurs des mots** d'une fenêtre de **contexte**.
- Calculer des **vecteurs de sous-mots** (vecteurs de caractères ou n-grams) et les assembler pour former un vecteur de substitution. [Bojanowski 2017]

$[?, ?, ?, \dots, ?] \rightarrow [0, 0, 0, \dots, 0]$

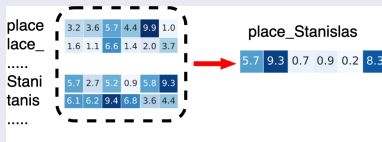
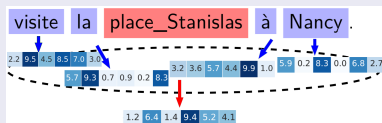


# Gestion des mots hors vocabulaire

## Différentes méthodes :

- Initialiser les vecteurs inconnus avec des **zéros** ou des valeurs **aléatoires**.
- Prendre le **centroïde** des **vecteurs des mots** d'une fenêtre de **contexte**.
- Calculer des **vecteurs de sous-mots** (vecteurs de caractères ou n-grams) et les assembler pour former un vecteur de substitution. [Bojanowski 2017]
- Exploiter des **données auxiliaires** (terminologies, thésaurus, bases de données lexicales telles que WordNet, etc.). [Pilehvar 2016]

$$[?, ?, ?, \dots, ?] \rightarrow [0, 0, 0, \dots, 0]$$

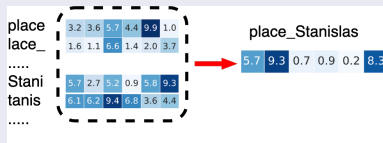
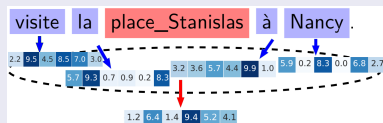


# Gestion des mots hors vocabulaire

## Différentes méthodes :

- Initialiser les vecteurs inconnus avec des **zéros** ou des valeurs aléatoires.
- Prendre le **centroïde des vecteurs des mots** d'une fenêtre de contexte.
- Calculer des **vecteurs de sous-mots** (vecteurs de caractères ou n-grams) et les assembler pour former un vecteur de substitution. [Bojanowski 2017]
- Exploiter des **données auxiliaires** (terminologies, thésaurus, bases de données lexicales telles que WordNet, etc.) [Pilehvar 2016]

$[?, ?, ?, \dots, ?] \rightarrow [0, 0, 0, \dots, 0]$



# Évaluation

## Évaluation des vecteurs de mots :

- **sur des tâches en aval** : classification de documents (dont analyse de sentiment), reconnaissance d'entités nommées, étiquetage morpho-syntaxique, traduction automatique, etc.
- **intrinsèque** : tâches de similarité et d'analogie.

→ Évaluation intrinsèque des méthodes de substitution.

# Évaluation

- **Corpus** : text8<sup>1</sup> pour l'entraînement et le calcul des substituts.
- **Paramètres** : par défaut de FastText et l'implémentation de Gensim pour Word2vec. Fenêtre de 2 pour calculer des centroïdes.

---

<sup>1</sup><http://mattmahoney.net/dc/textdata>

# Similarités original - substitut

## Objectif

Mesurer à quel point un vecteur construit (le substitut) est proche du vecteur réel du mot (original).

## Méthode

- Sélection aléatoire d'un ensemble de mots
- Pour chaque mot :
  - calcul du vecteur de substitution
  - similarité cosinus entre le vecteur de substitution et l'original
- Statistiques descriptives de la distribution des similarités cosinus



## Similarités original - substitut

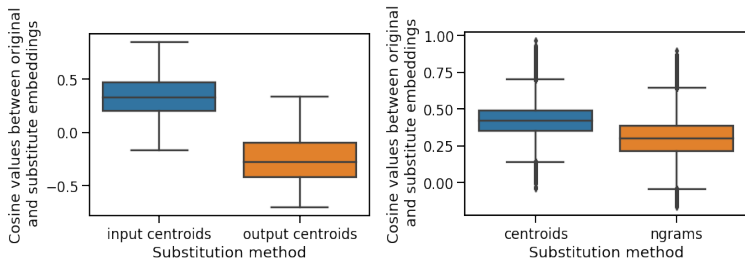


Figure 4: Similarité cosinus entre les vecteurs substitués et originaux pour Word2vec (gauche) et FastText (droite).

# Similarités paire d'originaux - paire de substituts

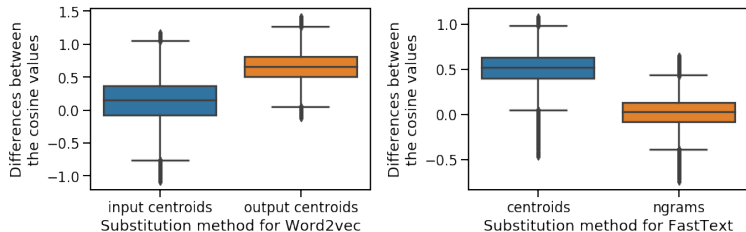
## Objectif

Mesurer à quel point deux mots dont les représentations vectorielles sont proches sont aussi proches dans l'espace des substituts.

## Méthode

- Sélection aléatoire d'un ensemble de paires de mots
- Pour chaque paire de mots :
  - calcul des vecteurs de substitution
  - calculer  $S_1$ , le cosinus entre les deux vecteurs originaux
  - calculer  $S_2$ , le cosinus entre les deux vecteurs de substitution
  - calculer  $S_1 - S_2$
- Statistiques descriptives de la distribution de  $S_1 - S_2$

## Similarités paire d'originaux - paire de substituts



**Figure 5:** Différence entre la similarité cosinus obtenue avec les vecteurs originaux et celle obtenue avec les substituts, pour Word2vec (gauche) et FastText (droite).

## Tâches de similarité

À chaque paire de mots est associé un score de similarité :

computer	keyboard	7.62
Jerusalem	Israel	8.46
planet	galaxy	8.11

### Évaluation

Score de **corrélation de rang** entre le classement donné par l'évaluation humaine et celui donné par les vecteurs.

## Tâches de similarité

**Table 1:** Scores de similarité quand les vecteurs de chaque paire sont issus du même espace : celui des vecteurs originaux ou celui des substituts.

classement 1	classement 2	wordsim	mturk	mc	rg
w2v_wc	humain	0.3501	0.1941	<b>0.3658</b>	<b>0.2654</b>
w2v_cc	humain	<b>0.3936</b>	<b>0.2236</b>	0.3562	0.2562
w2v_wc	w2v	<b>0.6551</b>	<b>0.4879</b>	<b>0.7018</b>	<b>0.6479</b>
w2v_cc	w2v	0.6212	0.4705	0.6155	0.5615
ft_ngrams	humain	0.3280	<b>0.4036</b>	<b>0.5000</b>	<b>0.4628</b>
ft_wc	humain	<b>0.3393</b>	0.2554	0.4760	0.4243
ft_ngrams	ft	0.4701	<b>0.4962</b>	0.5452	<b>0.5779</b>
ft_wc	ft	<b>0.5332</b>	0.4506	<b>0.7810</b>	0.5430

## Tâches de similarité

**Table 2:** Scores de similarité quand, pour chaque paire, les deux vecteurs sont l'original et le substitut.

classement 1	classement 2	wordsim	mturk	mc	rg
w2v_wc / w2v	humain	<b>0.3678</b>	<b>0.2904</b>	<b>0.2488</b>	<b>0.3740</b>
w2v_cc / w2v	humain	0.2652	0.2623	0.1640	0.0992
w2v_wc / w2v	w2v	<b>0.6009</b>	<b>0.5443</b>	<b>0.6654</b>	<b>0.5996</b>
w2v_cc / w2v	w2v	0.4855	0.4033	0.5167	0.1142
ft_wc / ft	humain	0.3568	0.3535	0.4613	0.4717
ft_ngrams / ft	humain	<b>0.4346</b>	<b>0.4673</b>	<b>0.5100</b>	<b>0.5089</b>
ft_wc / ft	ft	0.5919	0.6209	0.6275	0.7108
ft_ngrams / ft	ft	<b>0.6628</b>	<b>0.6977</b>	<b>0.8251</b>	<b>0.7613</b>

# Tâches d'analogie

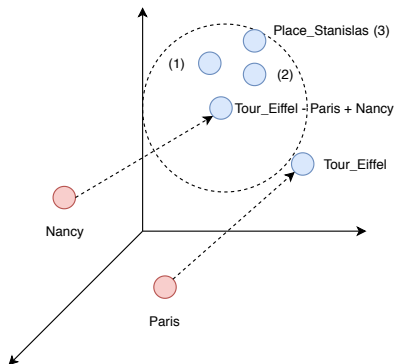


Figure 6: Analogies dans l'espace vectoriel.

Une tâche (sur 14) :  
ensemble de questions  
 $w_1 - w_2 + w_3 = ?$

Score par question :  
nombre de vecteurs  
plus proches de  
 $w_1 - w_2 + w_3$  que du  
vecteur cible.

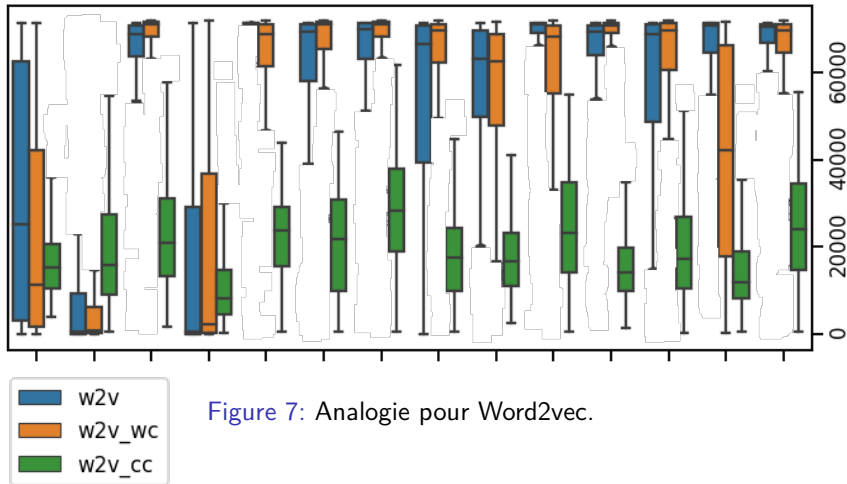


Figure 7: Analogie pour Word2vec.



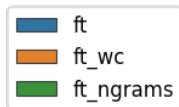
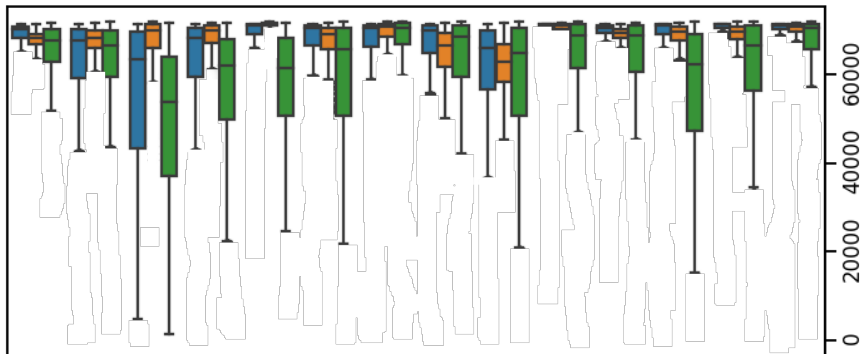


Figure 8: Analogie pour FastText.

## Conclusion

- Pas de méthode universellement meilleure.
- Résultats variables sur les différentes tâches :
  - Similarité : pour Word2vec, la matrice d'entrée donne de meilleurs résultats que la matrice de sortie.
  - Analogie :
    - pour Word2vec, la matrice de sortie donne des résultats avec moins de variabilité que la matrice d'entrée.
    - pour FastText, la méthode intégrée donne de meilleurs résultats au prix d'une plus grande variabilité.

## Perspectives

Utiliser des méthodes d'ensemble.

## Bibliographie

- Bojanowski P. et al. Enriching word vectors with sub-word information. *TACL* 2017.
- Mikolov T. et al. Distributed representations of words and phrases and their compositionality. *Advances in Neural Information Processing Systems* 2013.
- Pilehvar M. et al. Inducing embeddings for rare and unseen words by leveraging lexical resources. *ACL* 2017.
- Chu, C. et S. Kurohashi. Paraphrasing out-of-vocabulary words with word embeddings and semantic lexicons for low resource statistical machine translation. *LREC* 2016.

# Précisions pour FastText

## Deux méthodes de substitution :

- Centroïde des mots du contexte (comme pour Word2vec)
- Méthode intégrée : à partir des représentations de ngrams

## Entraînement

Les mots étudiés sont remplacés par le hash du mot :

- Une mise à jour des poids lorsque le mot est rencontré (représenté par le hash) ne met pas à jour les poids de ses sous-mots réels.
- Le hash du mot est suffisamment long pour que FastText puisse opérer (assez de sous-mots).