

Word embedding

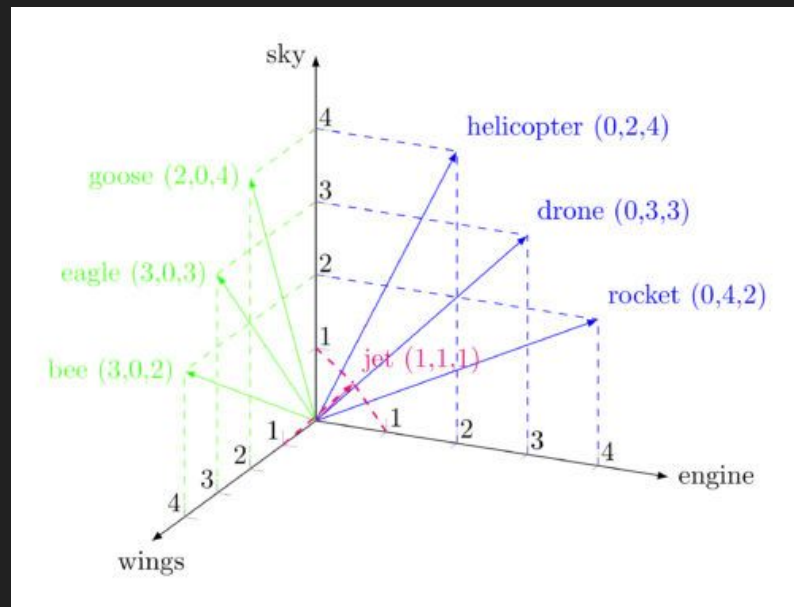
Traitement automatique du langage naturel

C'est quoi ?

Coder la sémantique des mots.

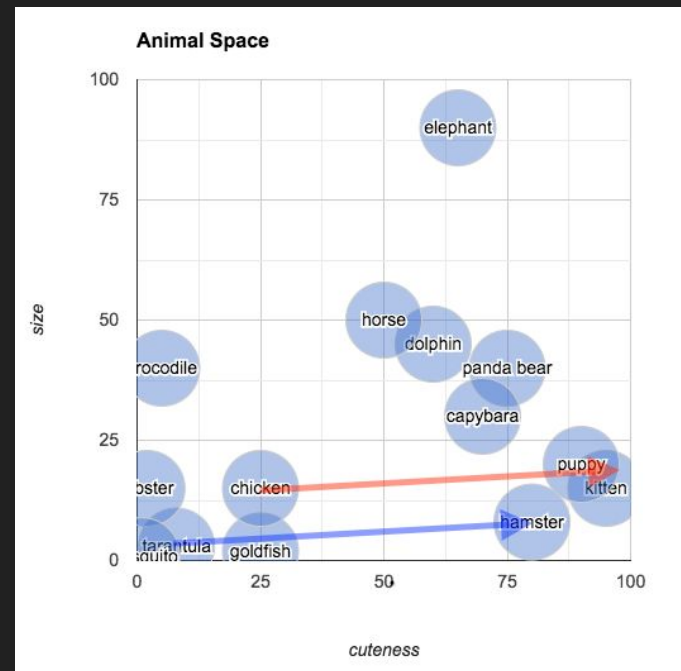
Moteurs de recherche, analyse de sentiment, similarité entre documents, analyse sémantique.

Apprentissage automatique.



Principe sur un exemple simple

	cuteness (0–100)	size (0–100)
kitten	95	15
hamster	80	8
tarantula	8	3
puppy	90	20
crocodile	5	40
dolphin	60	45
panda bear	75	40
lobster	2	15
capybara	70	30
elephant	65	90
mosquito	1	1
goldfish	25	2
horse	50	50
chicken	25	15



Sémantique distributionnelle

Pour un texte de taille
arbitraire ?

Un mot est défini par ses
voisins.

Il a fait bien **froid** hier.

Par contre, il va faire hyper **chaud**
aujourd'hui.

Il fera très **beau** demain !

Est-ce qu'il va faire très **froid** mardi ?

Sémantique distributionnelle

It was the best of times, it was the worst of times.

Réduction de la dimensionnalité.

Count-based.

g	DÉBUT _ was	it _ the	was _ best	the _ of	best _ times	of _ it	times _ was	was _ worst	worst _ times	of _ FIN
it	1	0	0	0	0	0	1	0	0	0
was	0	2	0	0	0	0	0	0	0	0
the	0	0	1	0	0	0	0	1	0	0
best	0	0	0	1	0	0	0	0	0	0
of	0	0	0	0	1	0	0	0	1	0
times	0	0	0	1	0	0	0	0	0	1
worst	0	0	0	1	0	0	0	0	0	0

Word2Vec

2013

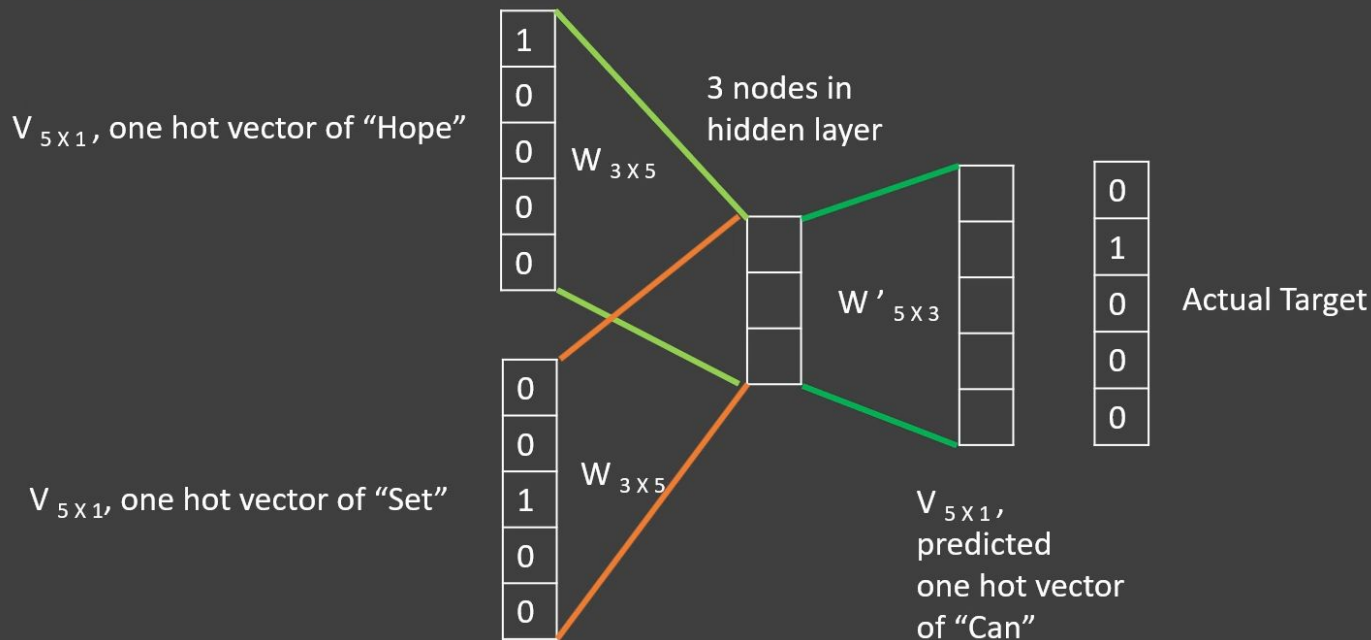
Prediction-based.

2 modèles légers et puissants : CBOW et Skip-gram.

Utiles pour l'apprentissage profond.

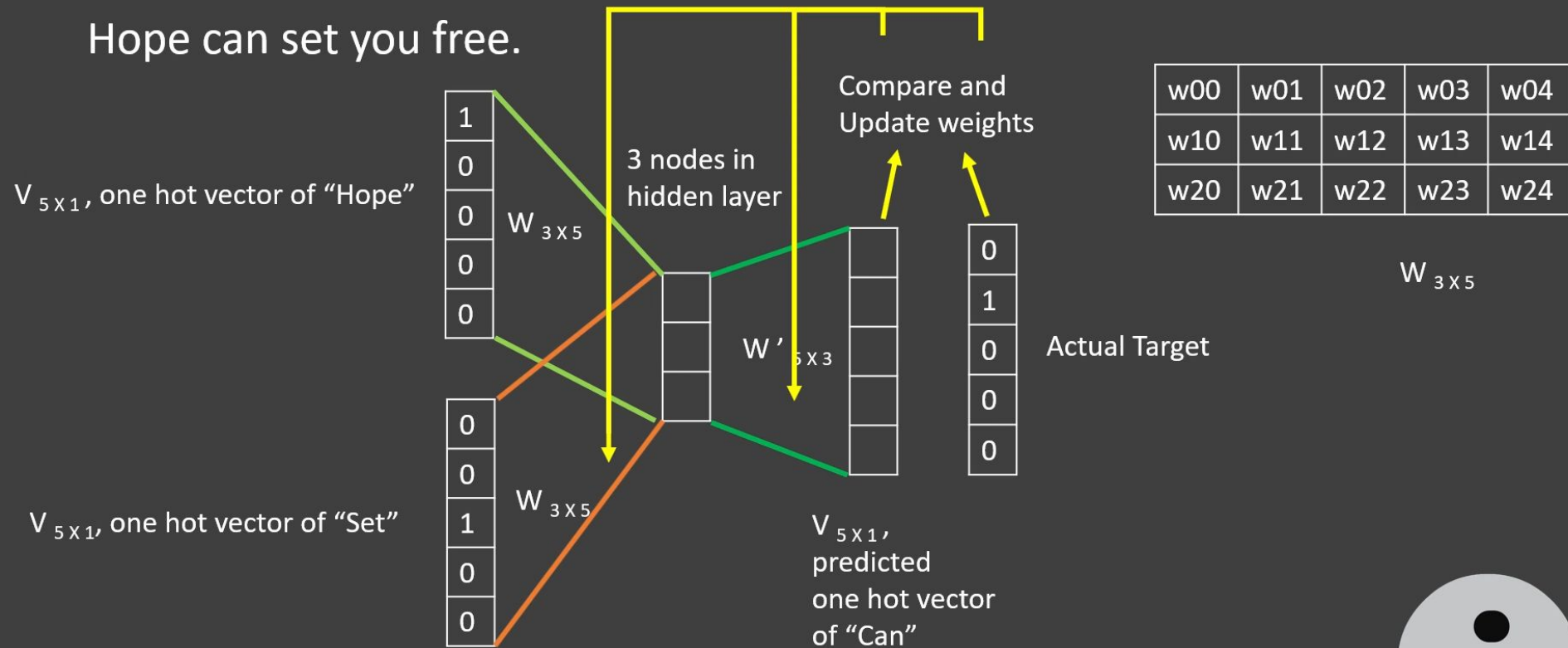
CBOW

Hope can set you free.



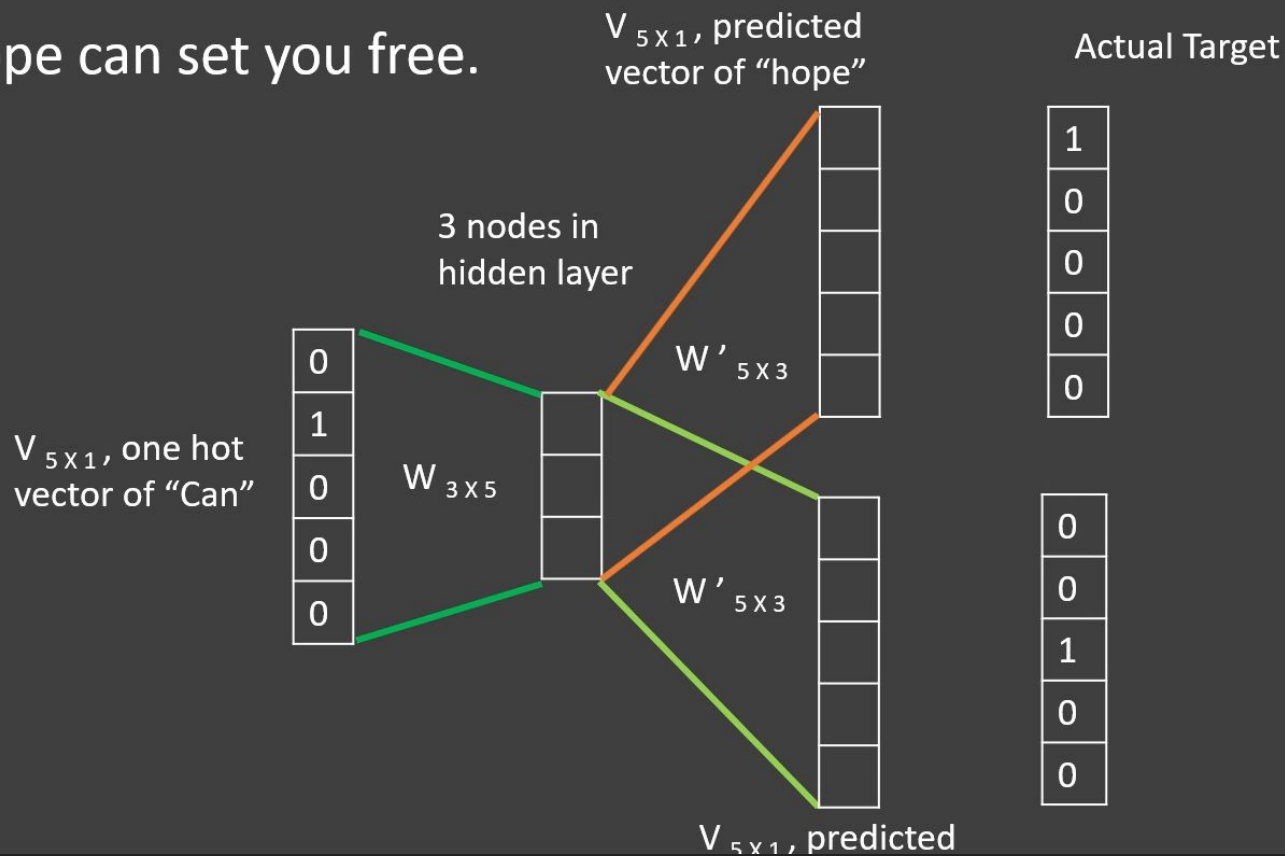
CBOW

Hope can set you free.



Skip-gram

Hope can set you free.



Sources

Efficient Estimation of Word Representations in Vector Space, Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean, 2013

<https://arxiv.org/abs/1301.3781>

Understanding word vectors, Allison Parish

<https://gist.github.com/aparrish/2f562e3737544cf29aaf1af30362f469>

https://en.wikipedia.org/wiki/Word_embedding, consulté le 31 mars 2024

https://en.wikipedia.org/wiki/Distributional_semantics, consulté le 31 mars 2024

<https://corpling.hypotheses.org/files/2018/04/3dplot-500x381.jpg> l'image de word embedding

<https://www.youtube.com/watch?v=UqRCEmrvlgQ>, The Semicolon (pour les images de CBOW et Skip-gram)