

NATURAL LANGUAGE PROCESSING

المعالجة اللغوية الطبيعية



المحتويات

				التطبيقات	العقبات و التحديات	تاريخ NLP	ما هو NLP	المحتويات	1) مقدمة
					البحث في النصوص	ملفات pdf	الملفات النصية	المكتبات	2) أساسيات NLP
T.Visualization	Syntactic Struc.	Matchers	Stopwords	NER	Stem & Lemm	POS	Sent. Segm.	Tokenization	3) أدوات NLP
	Dist. Similarity	Text Similarity	TF-IDF	BOW	Word2Vec	T. Vectors	Word embed	Word Meaning	4) المعالجة البسيطة
T. Generation	NGrams	Lexicons	GloVe	L. Modeling	NMF	LDA	T. Clustering	T. Classification	5) المعالجة المتقدمة
	Summarization & Snippets		Ans. Questions		Auto Correct	Vader	Naïve Bayes	Sent. Analysis	
Search Engine	Relative Extraction		Information Retrieval		Information Extraction		Data Scraping	Tweet Collecting	6) تجميع البيانات
					Rec NN\TNN	GRU	LSTM	Seq to Seq	7) RNN
Chat Bot	Gensim	FastText	Bert	Transformer	Attention Model	T. Forcing	CNN	Word Cloud	8) تكتيكات حديثة

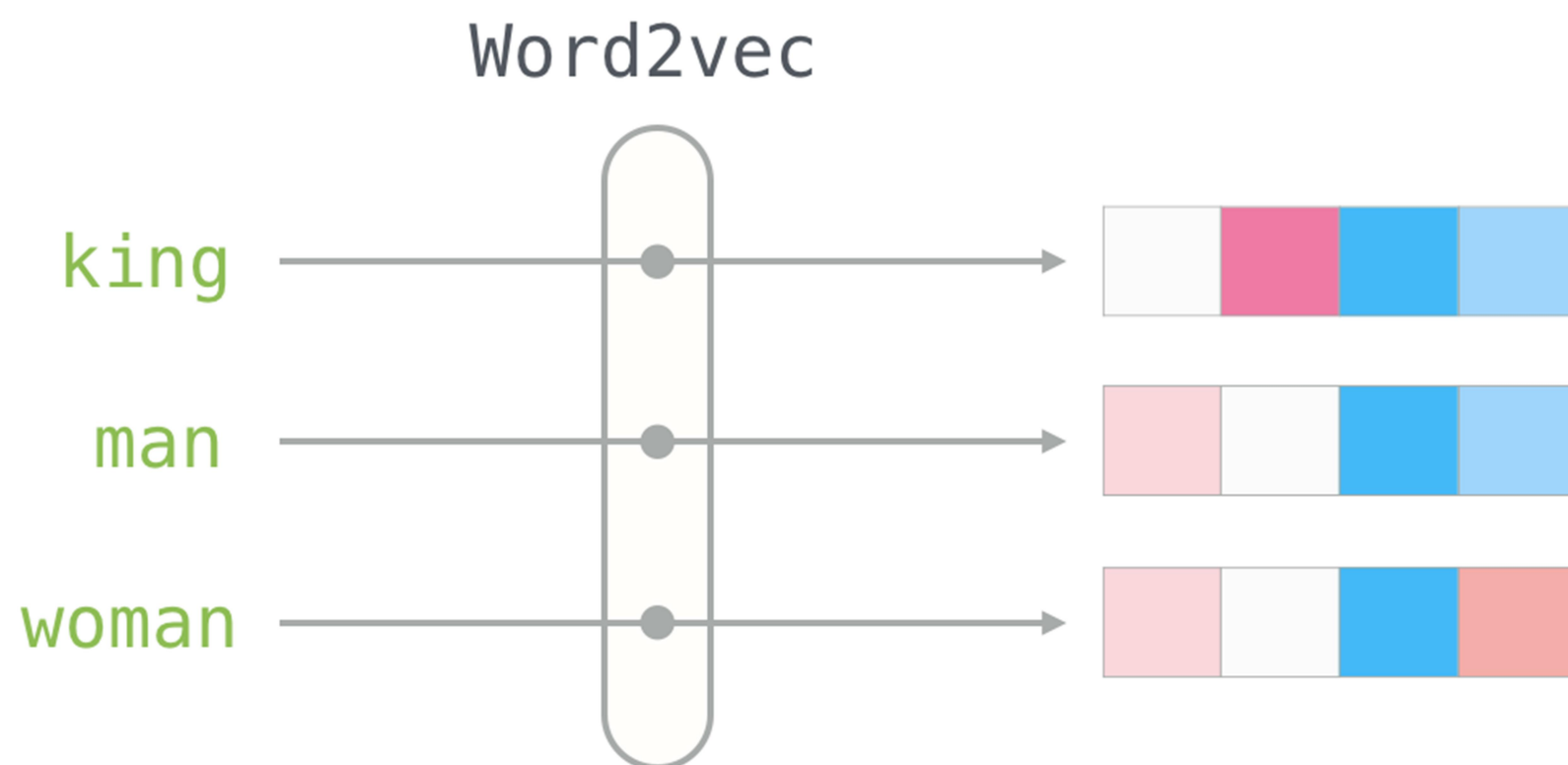
القسم الرابع : المعالجة البسيطة للنصوص

الجزء الرابع : Word2Vec

=====

اداة word2vec هي عبارة عن شبكة عصبية من طبقتين و التي تقوم بمعالجة النصوص .

يكون المدخل لها هو باقة الكلمات text corpus اما المخرج فهي كمية من المصفوفات الخاصة بالفيتشرز للنص



فهي ببساطة شبكة عصبية , يتم تدريبها علي أساس تضمين الكلمات , وهدفها حساب مدي أهمية و قيمة كل كلمة في الجملة ,
و من ثم , نقوم باستنتاج الكلمة الباقية

و المهمة الأساسية لأداة word2vec هي عمل تجميع grouping للمصفوفات للكلمات المتشابهة و المتماثلة و المرتبطة معا , وهو ما يتم عبر التشابهات الرياضية لكل كلمة

و هذه التشابهات و التناظرات تشبه (رجل - صبي) == (امرأة - فتاة) .

كذلك انها ستعرف ان هذه الكلمة مفردة و هذه جمع , وهو ما يسهل لاحقا عمل صياغة كاملة لنصوص و معرفة هل المفترض استخدام مفرد ام جمع و هكذا

و ايضا اداة word2vec حينما تاخذ كمية كبيرة من البيانات , فليديها المقدرة علي توقع معاني الكلمات , بناء علي موقعها و سياقها . .

* * * * *

و منها نوعان أساسيان له هما : CBOW , skip-gram

أولا : تكنيك CBOW

وهو اختصار continuous bag of words , وهو مزيج بين تكنيك bag of words و تكنيك NGram

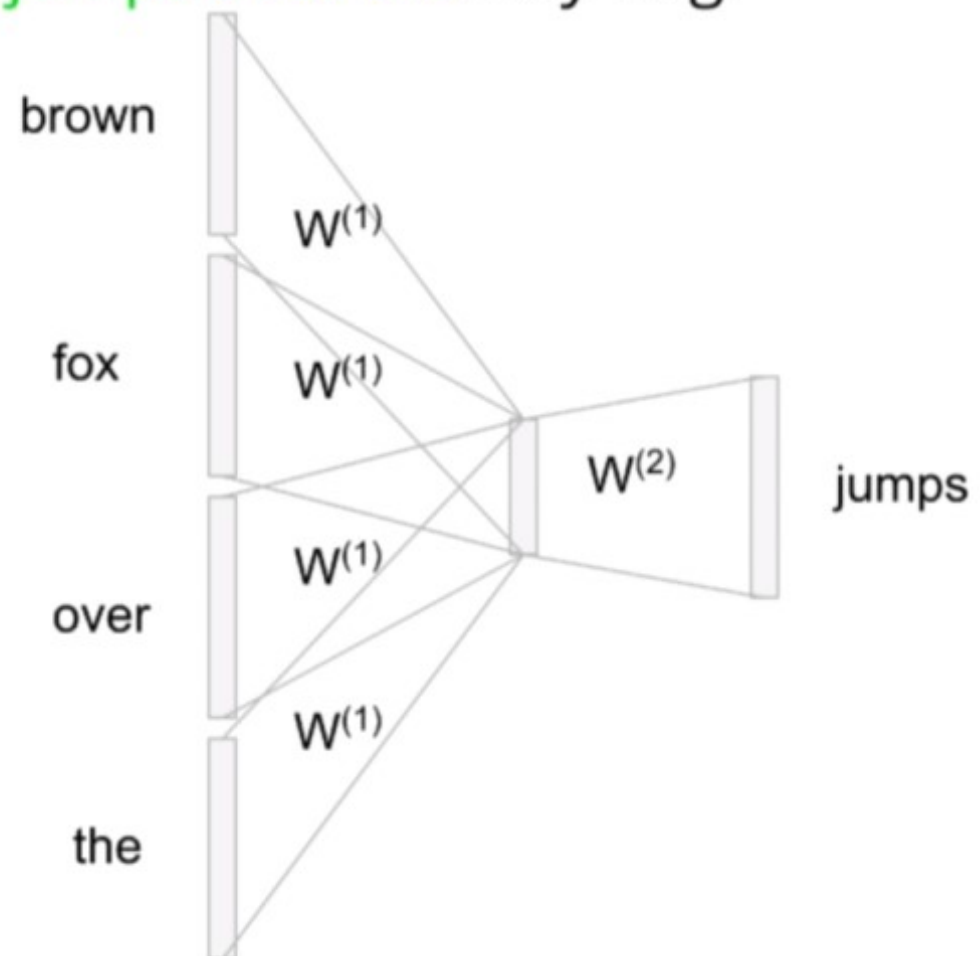
فكرة BOW تعتمد علي استخدام عدد من الكلمات الموجودة في النص و عمل ارقام 1 و 0 حسب تواجد كل كلمة , وفكرة NGram تعتمد علي استخدام كلمة او اكثر من الكلمات السابقة لاستنتاج الكلمة التالية , و بالتالي عبر استخدام CBOW يمكننا ان نستخدم اكثر من كلمة في نفس الجملة لاستنتاج كلمة محددة

و يمكن تعريفها انها شبكة عصبية , ولكن تستخدم لتوقع ما هي الكلمة الناقصة في جملة معينة (غالبا ما تكون الكلمة الأخيرة) , وذلك عبر حساب المصفوفات الضمنية embedding matrix للكلمات الداخلة , و منها يتم معالجتها للوصول للطبقة التالية في الشبكة , و اخيرا الطبقة النهائية بدالة تفعيل softmax لاختيار الكلمة الناقصة

و يتم اختيار عدد محدد من الكلمات التي سيتم استخدامها لاستنتاج الكلمة المطلوبة , فهنا يتم اختيار رقم 2 , اي كلمتين قبل الكلمة المطلوبة , وكلمتين بعدها , كما نري في المثال

CBOW - continuous bag of words

"The quick brown fox jumps over the lazy dog."



“Context size” could be considered 2 (or 4)

In practice, context size is usually set from 5-10 (on either side)

The input weight is $W^{(1)}$ for all input words (same weight used multiple times)

و يمكن زيادة الرقم اكثر من 2 , وغالبا ما يكون من 5 الي 10 , من كلا الاتجاهين

* * * * *

و من التطبيقات العملية لـ CBOW هي ما يسمى "تدوير المقال" article spinning

و الذي يقصد به ان نقوم بعمل تغيير في بعض كلمات المقالات , مع الاحتفاظ بالمعني , وذلك لعمل نوع من الاقتباس من بعض المواقع ووضعها في مواقع اخري , دون ان يقوم جوجل بتقليل تصنيف الموقع الثاني علي انه سارق للمحتوي

Text Before:

Content marketing is a strategic marketing approach focused on creating and distributing valuable, relevant, and consistent content to attract and retain a clearly defined audience — and, ultimately, to drive profitable customer action.

Instead of pitching your products or services, you are providing truly relevant and useful content to your prospects and customers to help them solve their issues.

Our annual research shows the vast majority of marketers are using content marketing. In fact, it is used by many prominent organizations in the world, including P&G, Microsoft, Cisco Systems, and John Deere. It's also developed and executed by small businesses and one-person shops around the globe. Why? Because it works.

Start

Text After:

Content advertising is a strategic advertising approach targeted on creating and distributing valuable, relevant, and steady content to appeal to and maintain a really described target audience — and, ultimately, to pressure profitable patron action.

Instead of pitching your merchandise or services, you are offering in reality applicable and useful content material to your prospects and clients to help them resolve their issues.

Our annual lookup suggests the vast majority of entrepreneurs are the usage of content marketing. In fact, it is used by way of many distinguished agencies in the world, together with P&G, Microsoft, Cisco Systems, and John Deere. It's also developed and carried out with the aid of small corporations and one-person shops round the globe. Why? Because it

Copy **Download**

Done

و من الأدوات التي تقوم بتطبيق article spinning هي الـ CBOW حيث اننا نقوم اولا بحذف بعض الكلمات من المقال , ثم نقوم هي باستنتاج الكلمة المطلوبة باستخدام الكلمات السابقة و التالية

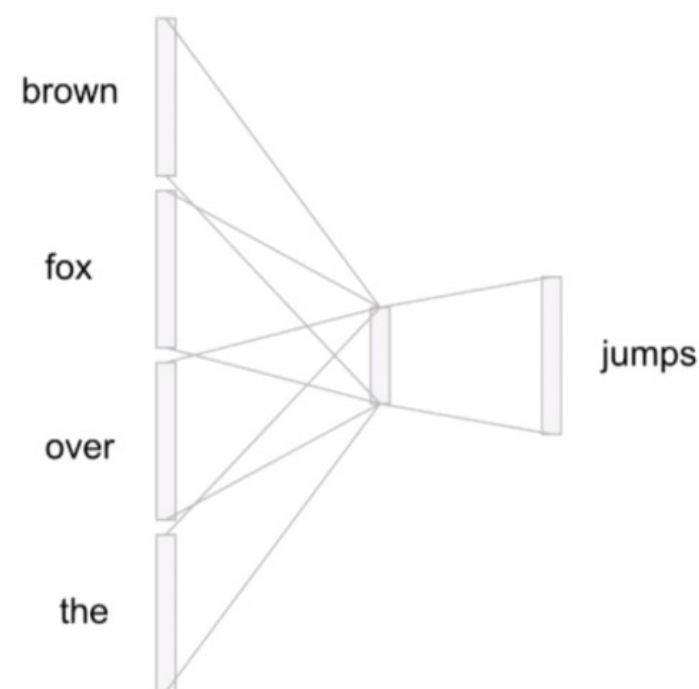
و تكون الفكرة الرياضية للـ CBOW اننا نقوم بما يشبه simple neural network و كان هناك عددا من المدخلات (ضعف الرقم الذي تم تحديده من الاتجاهين) , ثم نقوم بحساب قيم التضمين الخاصة بها, و ايجاد متوسط كل كلمة , ثم نقوم بعمل تدوير لهذه المصفوفة و ضربها في قيم التدريب h و باستخدام softmax يمكننا استنتاج الكلمة التالية

و عبر تدريب الموديل علي مئات الالاف من الكلمات , يمكن حساب قيم مناسبة للـ h ويمكن تطبيق الموديل بنجاح

$$h = \frac{1}{|C|} \sum_{c \in C} W^{(1)}_c$$

$$p(y | C) = \text{softmax}(W^{(2)T} h)$$

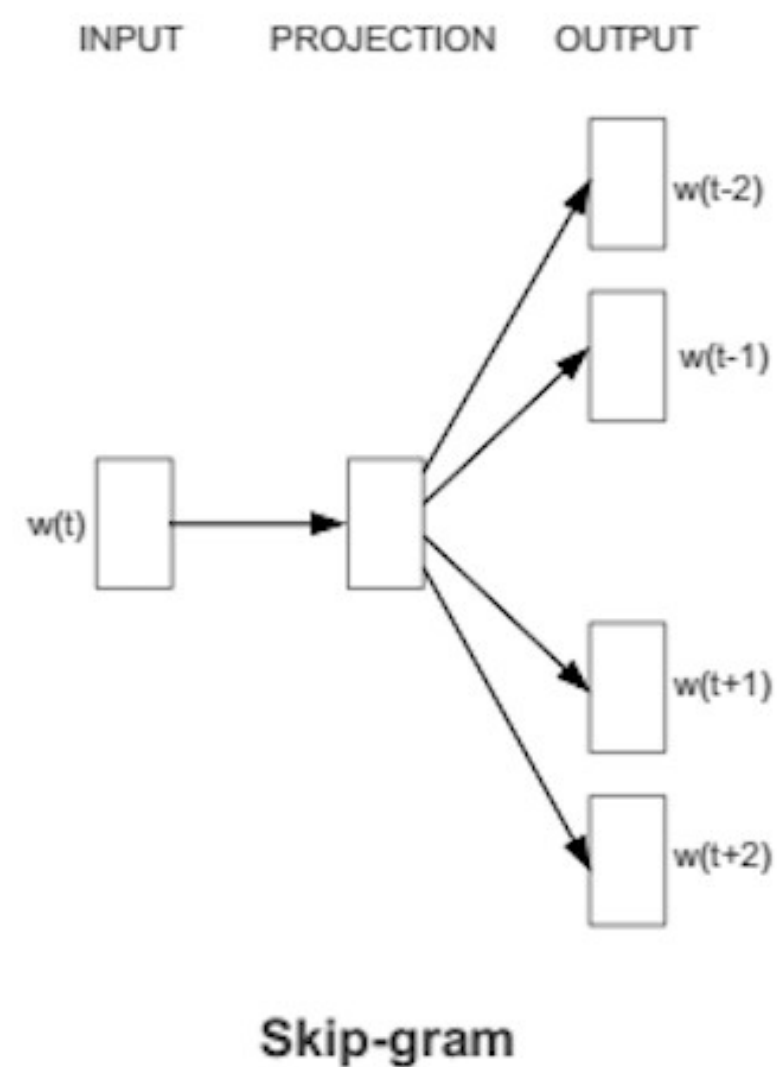
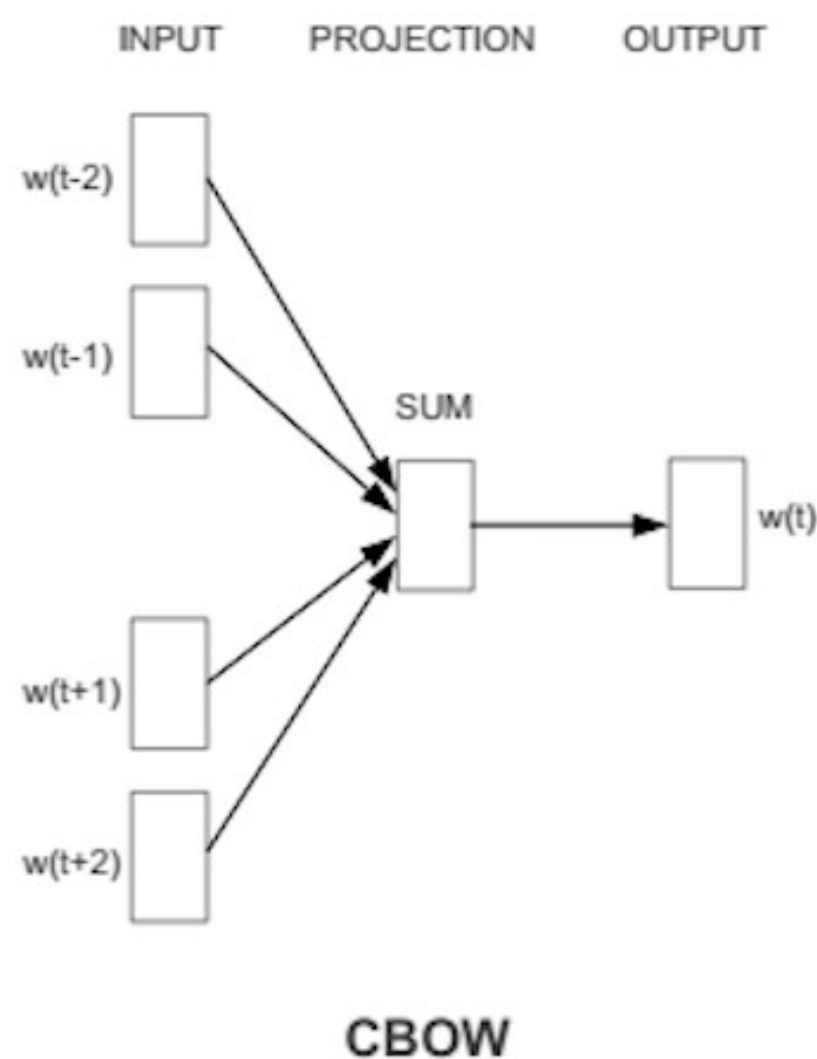
You could try to derive the gradients now, but we still have a few more modifications to make



* _ * _ * _ * _ * _ * _ * _ * _ * _ * _ * _ * _ * _

ثانيا : تكنيك SkipGram

إن كان الـ CBOW تقوم باستنتاج كلمة ناقصة عبر تناول كلمات الجملة فإن الـ skip-gram العكس , تستنتج جملة عبر كلمة



و تقوم فكرتها علي حساب العلاقة بين الكلمات و بعضها البعض , في المرة الأولى , تم حساب العلاقة بين الكلمة الاولى و الثانية this, is و بين الاولى و الثالثة a, this

و في المرة الثانية حينما كان التركيز علي كلمة is تم الحساب العلاقة بين الثانية و الاولى , و الثانية و الثالثة , و الثانية و الرابعة

Skip-gram for window 2					
	w+1	w+2			
This	is	a	NLP	Python	course
w-1		w+1	w+2		
This	is	a	NLP	Python	course
w-2	w-1		w+1	w+2	
This	is	a	NLP	Python	course
	w-2	w-1		w+1	w+2
This	is	a	NLP	Python	course
		w-2	w-1		w+1
This	is	a	NLP	Python	course
			w-2	w-1	
This	is	a	NLP	Python	course

(this, is) (this a)

(is, this) (is, a) (is, NLP)

(a, is) (a, this) (a, NLP) (a, Python)

(NLP, a) (NLP, is) (NLP, Python) (NLP, course)

(Python, NLP) (Python, a) (Python, course)

(course, Python) (course, NLP)

و هكذا , مع ملاحظة أنه يتم تحديد حد أقصى لعلاقة الكلمات مع بعضها البعض , فهنا الحد الاقصى هو 4

و تسمى هنا window 2 لأنه يتم وضع عدد 2 شباك امام و مثلهم في الخلف , و يتم تدريب الشبكة جيدا علي هذه الكلمات , حتي تتمكن من استنتاج الكلمات التالية عبر اعطائها كلمة واحدة

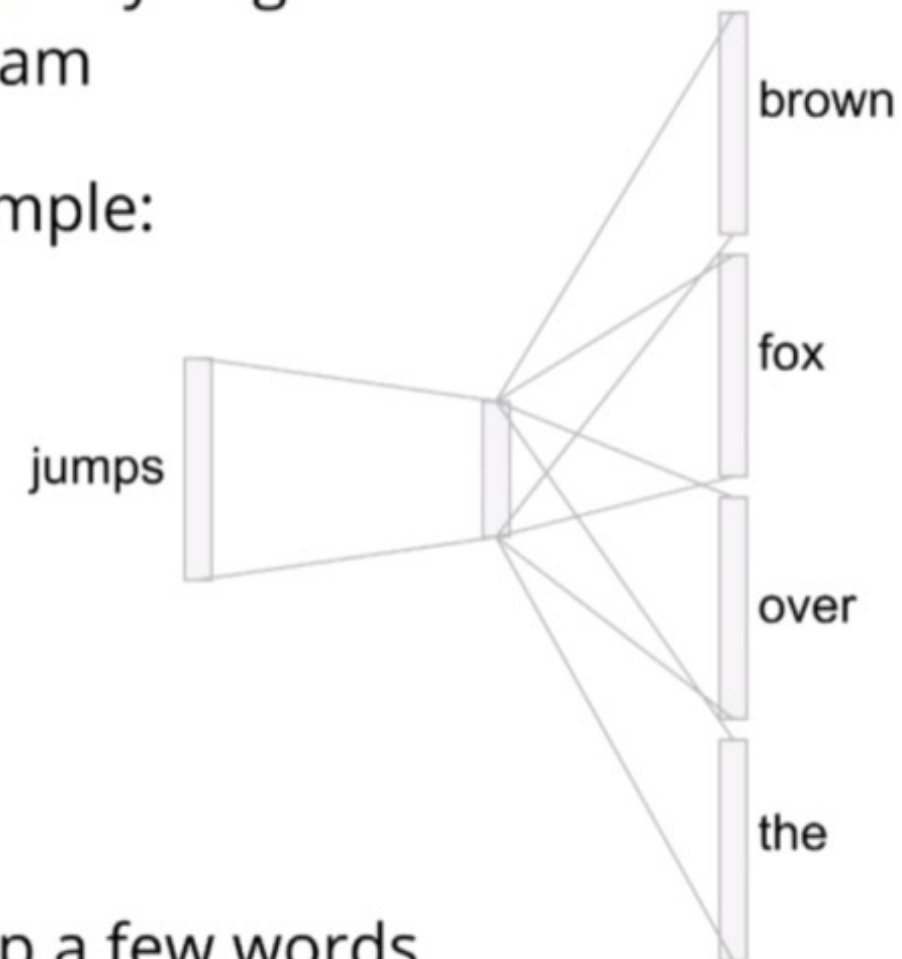
و ما يتم دخوله في الشبكة , هي قيم ال embedding لكل كلمة , ويتم حسابها بسهولة , عبر ايجاد ال onehotencoder لكل كلمة , وضربها في embedding matrix العامة

$$[0 \ 0 \ 0 \ 1 \ 0] \times \begin{bmatrix} 0.5 & 0.8 & 1.3 \\ 2.1 & 1.2 & 0.2 \\ 0.4 & 0.7 & 1.1 \\ 2.8 & 1.4 & 0.9 \\ 0.8 & 1.2 & 0.4 \end{bmatrix} = [2.8 \ 1.4 \ 0.9]$$

* * * * *

و تأتي الفكرة أصلا من فكرة ال bigram ففيها يمكن استخدام كلمة لاستنتاج كلمة تالية لها , اي من كلمة jumps يمكن استنتاج كلمة over , و لكن مع توسيع الفكرة يمكن عمل skipping اي تخطي لعدد من الكلمات بحيث نستنتج الكلمات الثالثة او الرابعة او السابقة , لذا تسمى skipgram

- "The quick brown fox jumps over the lazy dog."
- Helpful to think of it in terms of bigram
- Bigram model gives us 1 training sample:
jumps → over
- Skipgram gives us 3 additional training samples:
jumps → brown
jumps → fox
jumps → the
- Skipgram: *like* bigram, except we skip a few words



Source Text

Training Samples

The quick brown fox jumps over the lazy dog. ➡

(the, quick)
(the, brown)

The	quick	brown	fox	jumps over the lazy dog. ➡
-----	-------	-------	-----	----------------------------

(quick, the)
(quick, brown)
(quick, fox)

The	quick	brown	fox	jumps	over the lazy dog.	➡
-----	-------	-------	-----	-------	--------------------	---

(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

The quick brown fox jumps over the lazy dog. →

(fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)

و لكن ما هي الآلية التي سيتم اختيار الكلمات علي اساسها , كيف يمكن عبر معالجة كلمة ان نستنتج كلمة تالية لها ؟

هناك طريقتين , طريقة أقل نجاحا و هي الاختيار الهرمي , وطريقة أكثر نجاحا و هي : العينة السلبية

* * * * *

الاختيار الهرمي Hierarchical Softmax

و هي الطريقة الاولى لاختيار الكلمة المناسبة

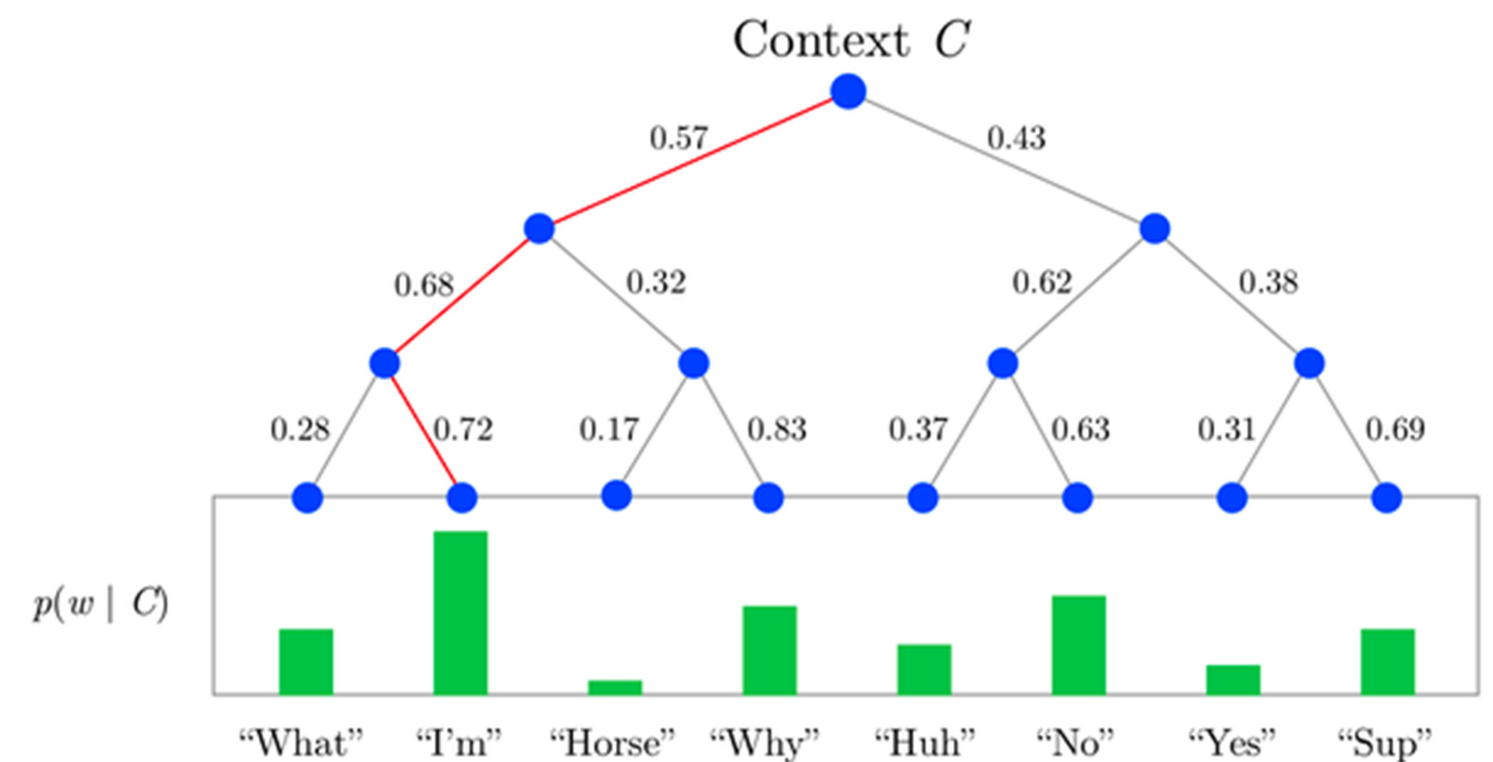
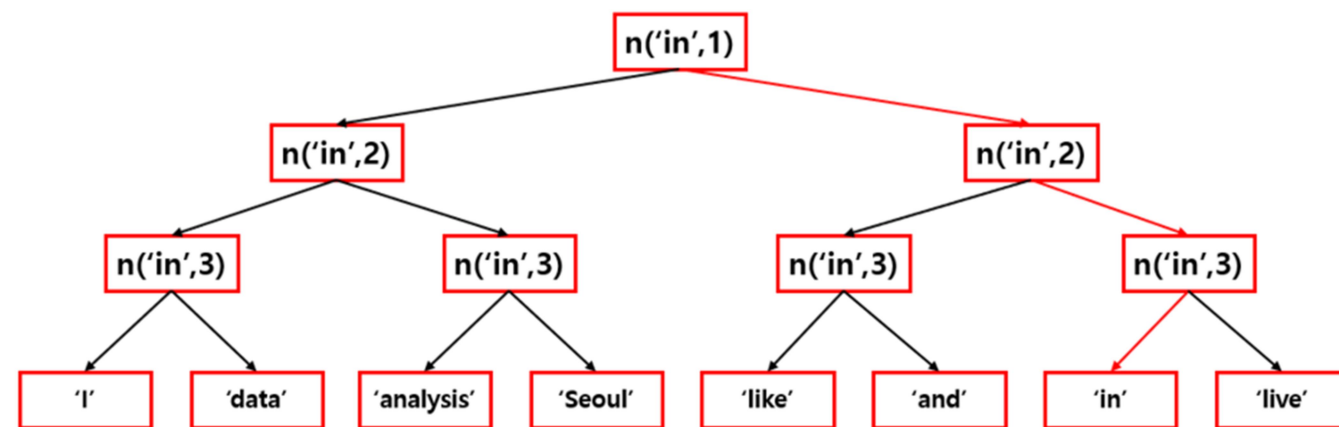
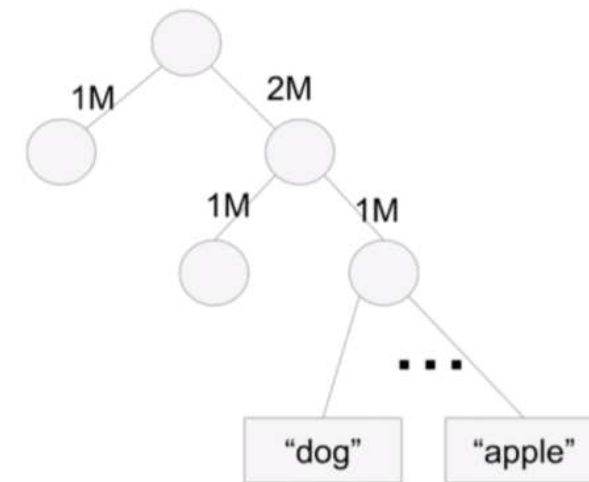
المعضلة تكمن في أن دالة ال softmax ستقوم بالاختيار عبر عدد ضخم من الكلمات , ففي جلوف هناك 400 ألف كلمة , وفي مكتبة word2vec هناك 3 مليون , وهذا معناه انه اثناء التصنيف ستكون هناك صعوبة جمة في اختيار الكلمة المناسبة , وستكون الدقة دائما قليلة

و السبب الاساسي ان دالة السوفت ماكس تعتمد علي مجموعة لوغاريتمات جميع الاصناف , فاذا قمنا بجمع لوغاريتمات 3 مليون كلمة , فسيكون المقام قيمة ضخمة جدا

- Calculating softmax is $O(V)$

$$p(y = j \mid x) = \frac{\exp(w_j^T x)}{\sum_{k=1}^V \exp(w_k^T x)}$$

لذا فإن فكرة الاختيار الهرمي قائمة علي تقسيم جميع الكلمات الموجودة لدينا علي شكل تقسيم هرمي كامل , بحيث يكون هناك مستويات في الكلمات , بحيث يكون هناك كلمات في المستويات العليا , وأخري في المستويات الاقل و هكذا



الطريقة الثانية لاختيار الكلمات , وهي الأكثر سرعة و دقة , وتسمى العينة السلبية negative sampling

و تأتي هذه الفكرة لتجنب عيوب الطريقة السابقة , حيث مشكلة التدريب علي مئات الالاف من الكلمات , وصعوبة الحصول علي كفاءة عالية

و الفكرة تقوم علي أنه بدلا من تدريب الموديل بحيث يكون هناك multi-classification بين الكلمة الصحيحة و مئات الالاف من الكلمات الخاطئة , فإننا نقوم بتحديد أولا بتحديد الكلمة الداخلة و هي jumps ثم تحديد الكلمات الصحيحة و هي (over , the , brown , fox) , وبدلا من تدريبها علي باقي الكلمات الأخرى في اللغة الانجليزية , فإننا نقول باختيار عدد من الكلمات الأخرى بشكل عشوائي , فيتم مثلا اختيار اربع كلمات اخري غير صحيحة بشكل عشوائي و لتكن (apple , Tokyo , boat , orange ,) , ثم تدريب الموديل بطريقة حساب احتمالية وجود الكلمات الصحيحة و زيادتها , ثم حساب احتمالية الكلمات الخطأ و تقليلها

و تكون معادلة الخطأ هي مجموعة لوغاريتمات احتمالية الكلمات الصحيحة , و لوغاريتمات 1 ناقص احتمالية الكلمات الخطأ

لا تنس ان لوغاريتم الـ 1 يساوي صفر , فمعادلة الخطأ ستؤول للصفر حينما تكون احتمالية الكلمات الصحيحة كبيرة و الخاطئة قليلة , وهنا نتمكن من تدريب الموديل علي اختيار الكلمات الصحيحة و الابتعاد عن الخاطئة

و بالنسبة لعدد الكلمات , غالبا ما يتم اختيار من 5 الي 10 كلمات

Input word: jumps

Target words: brown, fox, over, the

Negative samples: apple, orange, boat, tokyo

$$J = \log p(\text{brown} | \text{jumps}) + \log p(\text{fox} | \text{jumps}) + \\ \log p(\text{over} | \text{jumps}) + \log p(\text{the} | \text{jumps}) + \\ \log[1 - p(\text{apple} | \text{jumps})] + \log[1 - p(\text{orange} | \text{jumps})] + \\ \log[1 - p(\text{boat} | \text{jumps})] + \log[1 - p(\text{tokyo} | \text{jumps})]$$

و تكون المعادلة العامة :

$$J = \sum_{c \in C} \log \sigma(W^{(2)}_c^T W^{(1)}_{in}) + \sum_{n \in N} \log[1 - \sigma(W^{(2)}_n^T W^{(1)}_{in})]$$

و هناك فكرة تستخدم أحيانا , فلو كان لدينا كلمة مدخل وهي jumps و لدينا 4 كلمات صحيحة , فبدلاً من البحث عن عدد من الكلمات الخاطئة , فيمكن استبدال كلمة المدخل بكلمة واحدة خطأ , وليكن lighthouse , ثم تكون المعادلات الصحيحة هي العلاقة بين كلمة المدخل و الأربع كلمات الصحيحة , وتكون المعادلات الخطأ هي العلاقة بين كلمة المدخل الخاطئة و الأربع كلمات الصحيحة

The quick brown fox jumps over the lazy dog.

The quick brown fox lighthouse over the lazy dog.

+ve samples: jumps → brown, jumps → fox, ...

-ve samples: lighthouse → brown, lighthouse → fox, ...

و حينما نقوم باختيار الكلمات الخاطئة , من الممكن ان يتم هذا بشكل عشوائي , او أن نتبع معايير معينة بحيث تكون الكلمات الخاطئة ليست نادرة تماماً , حتي يتمكن الموديل من استبعاد الكلمات المنتشرة

و غالبا ما يتم استخدام هذا القانون , والذي يعتمد علي نسبة تواجد الكلمة اصلا في النص , فلو كانت كلمة نادرة , فستكون قيمة $\tilde{p}(w)$ برقم قليل , وبالتالي قيمة Pdrop ستكون قليلة تقترب من الصفر , والعكس صحيح

$$p_{drop}(w) = 1 - \sqrt{\frac{threshold}{\tilde{p}(w)}}$$

- Typical threshold = 10^{-5}
- Ex. if $p(w) = 10^{-5}$, then $p_{drop}(w) = 1 - 1 = 0$
- Ex. if $p(w) = 0.1$, then $p_{drop}(w) = 1 - 10^{-2} = 0.99$

حيث قيمة $\tilde{p}(w)$ تساوي

$$\tilde{p}(w) = \frac{count(w)^{0.75}}{\sum_{w'} count(w')^{0.75}}$$

*_*_*_*_*_*_*_*_*_*_*_*_*_*_*_*_*_*_*_*_*_**

و نري مثالين عمليين مبسطين لفكرة CBOW & Skip-gram في ملفي 4.4.1 & 4.4.2
بالإضافة إلي أن word2vec سيتم استخدامها بشكل افضل مع شرح مكتبة gensim