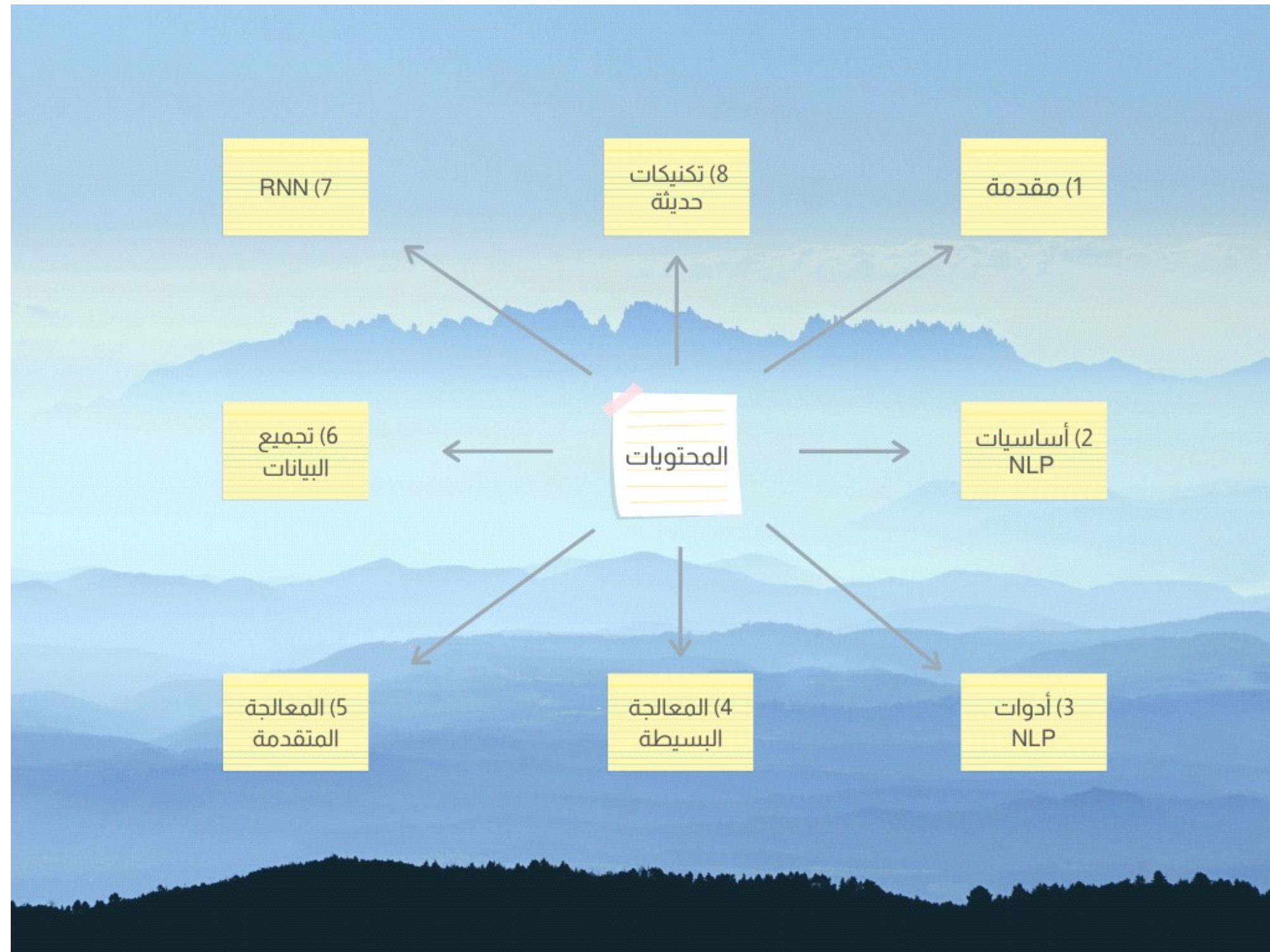


# NATURAL LANGUAGE PROCESSING

# 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) تكنيكات حديثة

## القسم الأول : مقدمة

### الجزء الأول : المحتويات

و هنا عدد من المفاهيم الهامة , التي يجب أن نراجع عليها سريعا :

#### • تجميع البيانات النصية **collecting text data**

و هي التي تتم عبر قراءة الملفات النصية , سواء كانت بامتداد `txt` , `pdf` , `csv` , `xlsx` و اثناء القراءة يمكن ان نقوم بعمليات بحث داخلها , واستخلاص بيانات محددة و ليست كلها

#### • معالجة النصوص

و هي التي تتم عبر مكتبات عديدة مثل `re` , `spacy` , `nltk` و التي تقوم بجميع عمليات البحث و المعالجة المطلوبة

## • جذور الكلمة stemming & lemmatization

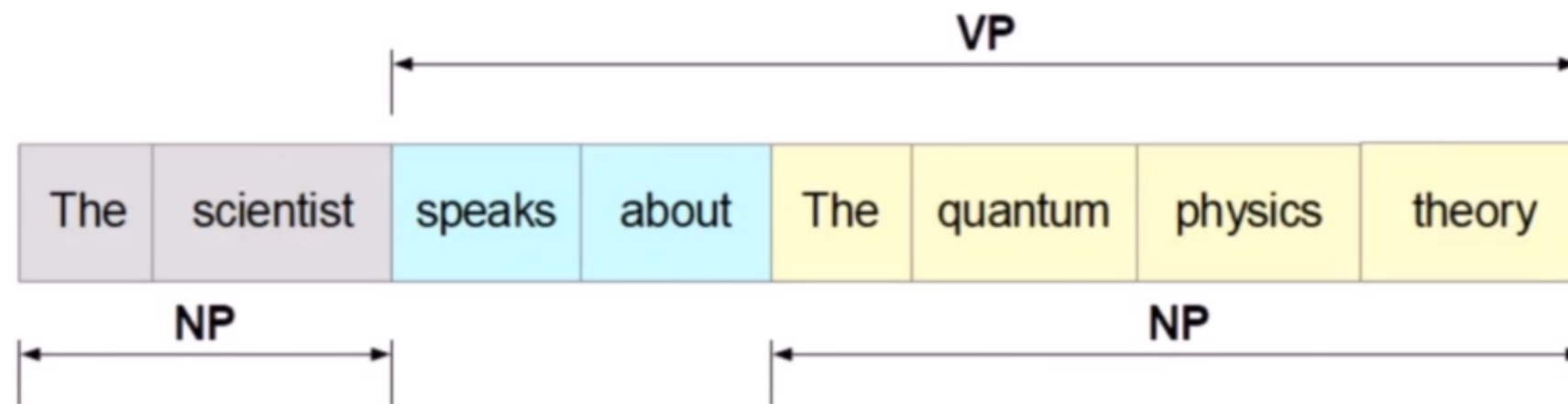
و هي خطوات ايجاد اصل الكلمة و جذرها , و حذف كل الاضافات عليها , و ايجاد الفعل الاصلي لها

## • جزء النصوص POS

و هي تحديد نوع هذه الكلمة في قواعد اللغة , هل هي اسم او فعل او صفة , و ربطها بباقي الجمل

## ● تقطيع الجمل chunking

و هي تحديد اي كلمات تنتمي لبعضها البعض , و كيف تقوم بتكوين الجمل معا



**NP** : Noun Phrase

**VP** : Verb Phrase



## • جعبة الكلمات Bag of Words

و هي طريقة جمع كل الكلمات الموجودة في النص , ثم تحديد هل الجملة الاولى او القطعة الاولى فيها هذه الكلمات ام لا

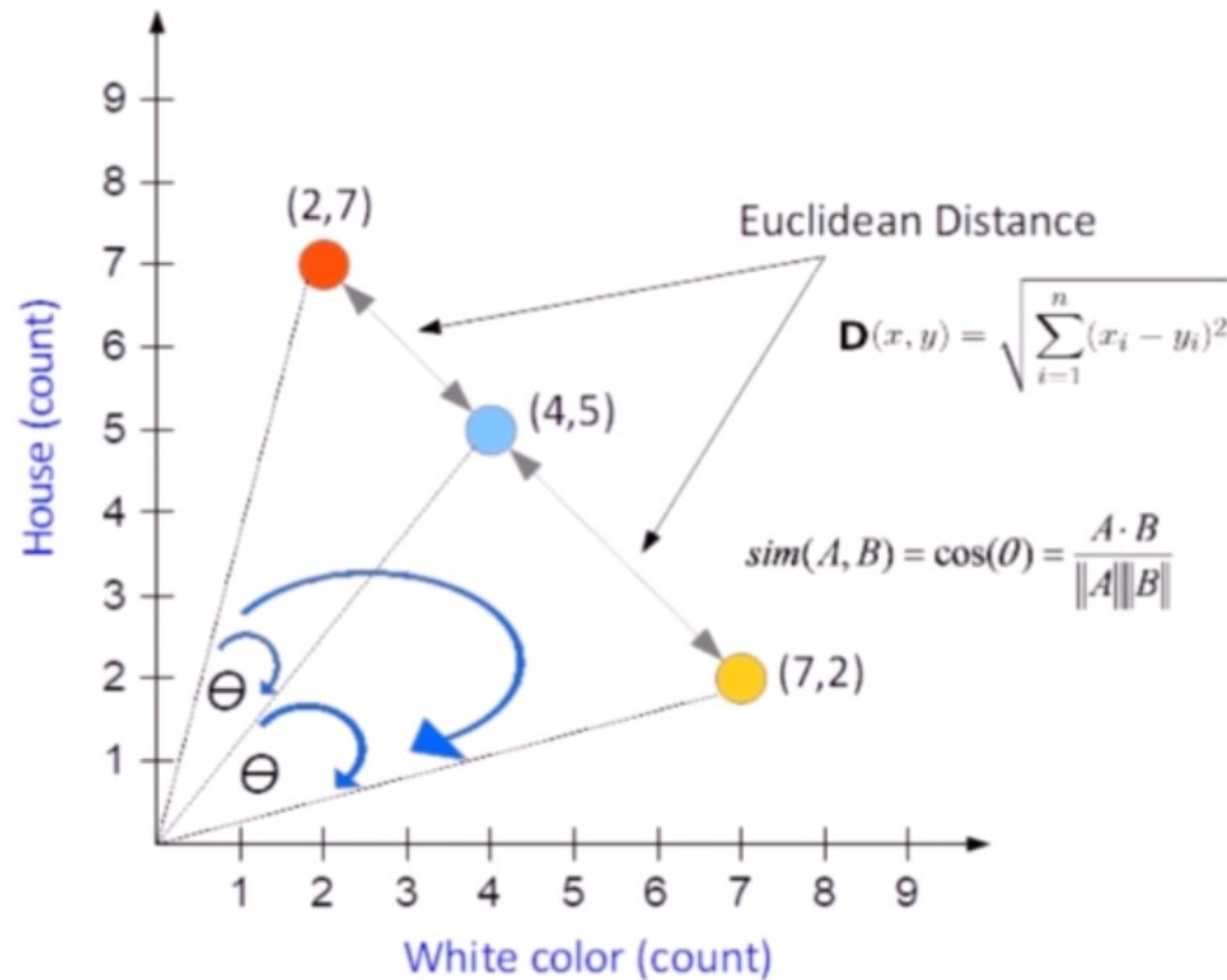
Term Document Matrix (TDM)			
Terms	Text 1	Text 2	Text 3
An	1	0	0
epic	1	0	0
adventure	1	0	0
in	1	1	0
time	1	1	0
space	1	1	1
and	1	1	1
life	1	0	0
the	0	1	0
cosmos	0	1	0
is	0	1	0
infinite	0	1	0
I	0	0	1
like	0	0	1
books	0	0	1
about	0	0	1
history	0	0	1

Document vector

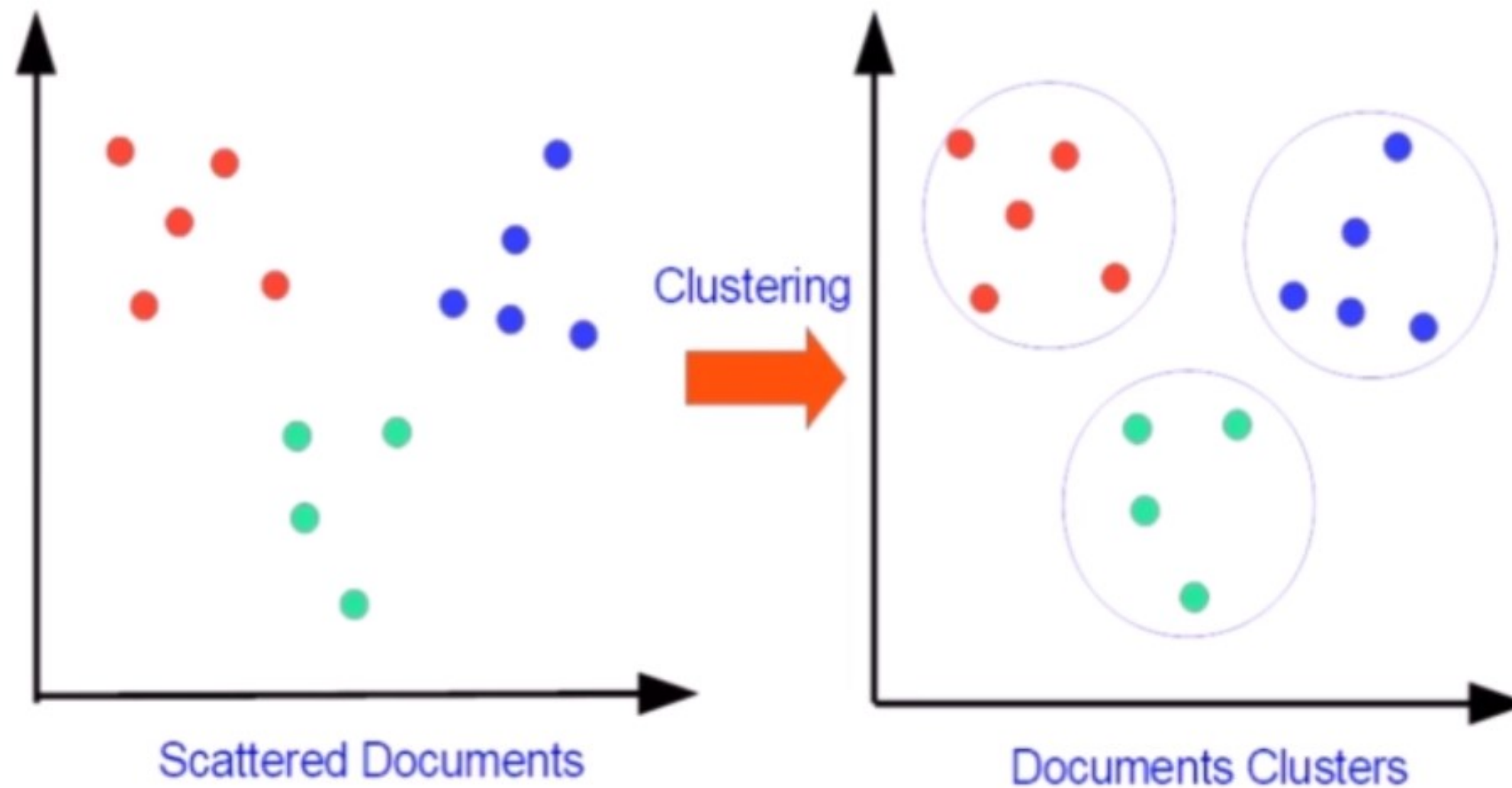
Word vector

## • تشابه الكلمات text similarity

و هي تشير الي مدي اقتراب معاني الكلمة الاولى من الثانية , عبر المعاني المضمنة embedding , ويتم حسابها باكثر من طريقة منها cos angle , او المسافة المطلقة بفيثاغورث



- **عناقيد الكلمات :**  
وهي استخدام تعليم بدون اشراف لتصنيف الكلمات الي عناقيد , بناء علي معانيها المضمنة





## ● تكنيك LDA :

و هي احد اساليب التعليم بدون اشراف , حيث نقوم بتقسيم الجمل الي عدد من الاقسام , بناء علي المعاني الدفينة في الجمل

## ● تكنيك NGrams :

● و هو اسلوب التعامل مع عدد معين من الكلمات السابقة و الكلمات التالية , لاستنتاج الكلمة الجديدة , او للتصنيف

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)

- تضمين الكلمات : وهو اظهار ال vectors الخاصة بكل كلمة , والتي يمكن من خلالها معرفة مدي اقتراب او ابتعاد الكلمات

