

Comunidad MyFuture-AI

# Introducción al Procesamiento de Lenguaje Natural (NLP)

By Laura López

# Agenda

Esto es lo que abarcaremos:

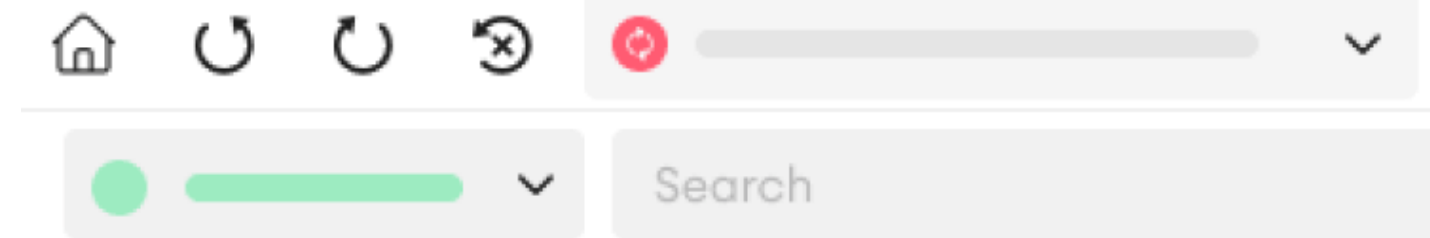
- Embeddings
- Preprocesamiento de texto
- Análisis NLP
- Práctica

# Procesamiento de Lenguaje Natural (NLP)



# Procesamiento de Lenguaje Natural (NLP)

“Es una tecnología de machine learning que brinda a las computadoras la capacidad de interpretar, manipular y comprender el lenguaje humano.” - Amazon



Despite the rapid advances in AI, computers are still challenged in matching the precision of human language. The training data here is as important as the accurate the input data annotation, model prediction.

How do we annotate data, though? To go with this one, but it all depends on the purposes of this article, we'll take a look at text boxes as one of the most extensively used techniques. Moving forward, we'll work on the following:

# ¿Cómo representas un texto con números?

Asignando a cada carácter un valor

Caracteres ASCII de control			Caracteres ASCII imprimibles			ASCII extendido (Página de código 437)										
00	NULL	(carácter nulo)	32	espacio	64	@	96	`	128	Ç	160	á	192	Ł	224	Ó
01	SOH	(inicio encabezado)	33	!	65	A	97	a	129	ü	161	í	193	ł	225	ô
02	STX	(inicio texto)	34	"	66	B	98	b	130	ë	162	ó	194	Ł	226	ö
03	ETX	(fin de texto)	35	#	67	C	99	c	131	â	163	ú	195	ł	227	õ
04	EOT	(fin transmisión)	36	\$	68	D	100	d	132	ä	164	ñ	196	—	228	ö
05	ENQ	(consulta)	37	%	69	E	101	e	133	à	165	Ñ	197	Ł	229	ő
06	ACK	(reconocimiento)	38	&	70	F	102	f	134	â	166	ª	198	ä	230	µ
07	BEL	(timbre)	39	'	71	G	103	g	135	ç	167	º	199	Å	231	þ
08	BS	(retroceso)	40	(	72	H	104	h	136	ê	168	¿	200	Ł	232	ƒ
09	HT	(tab horizontal)	41	)	73	I	105	i	137	ë	169	®	201	Ł	233	Ů
10	LF	(nueva línea)	42	*	74	J	106	j	138	è	170	¬	202	Ł	234	Ů
11	VT	(tab vertical)	43	+	75	K	107	k	139	ì	171	½	203	Ł	235	Ů
12	FF	(nueva página)	44	,	76	L	108	l	140	í	172	¾	204	Ł	236	ý
13	CR	(retorno de carro)	45	-	77	M	109	m	141	î	173	ı	205	Ł	237	ÿ
14	SO	(desplaza afuera)	46	.	78	N	110	n	142	Ä	174	«	206	Ł	238	—
15	SI	(desplaza adentro)	47	/	79	O	111	o	143	Å	175	»	207	Ł	239	·
16	DLE	(esc.vínculo datos)	48	0	80	P	112	p	144	É	176	÷	208	Ł	240	≡
17	DC1	(control disp. 1)	49	1	81	Q	113	q	145	æ	177	×	209	Ł	241	±
18	DC2	(control disp. 2)	50	2	82	R	114	r	146	Æ	178	÷	210	Ł	242	≡
19	DC3	(control disp. 3)	51	3	83	S	115	s	147	ö	179	÷	211	Ł	243	¾
20	DC4	(control disp. 4)	52	4	84	T	116	t	148	ö	180	÷	212	Ł	244	¶
21	NAK	(conf. negativa)	53	5	85	U	117	u	149	ò	181	À	213	Ł	245	§
22	SYN	(inactividad sinc)	54	6	86	V	118	v	150	û	182	Á	214	Ł	246	÷
23	ETB	(fin bloque trans)	55	7	87	W	119	w	151	ù	183	Â	215	Ł	247	º
24	CAN	(cancelar)	56	8	88	X	120	x	152	ÿ	184	©	216	Ł	248	º
25	EM	(fin del medio)	57	9	89	Y	121	y	153	Ö	185	÷	217	Ł	249	·
26	SUB	(sustitución)	58	:	90	Z	122	z	154	Û	186	÷	218	Ł	250	·
27	ESC	(escape)	59	;	91	[	123	{	155	ø	187	÷	219	Ł	251	·

Tomada de: <https://elcodigoascii.com.ar>

Tiene sentido, pero ¿qué pasa si tengo un documento?

¿Cómo representas  
un texto con  
números?

Asignando a cada palabra un  
valor o posición de un vector

**My Future IA**

11

10

23

Pero ¿**My** es mayor que **Future**?

# ¿Cómo representas un texto con números?

**Vectores para representar las palabras y sentido semántico**

**One-Hot encoding**

[My, Future, IA]

My: [1,0,0]

Future:[0,1,0]

IA:[0,0,1]

Todas las palabras tendrían la misma distancia, ¿qué pasa si tengo palabras con un contexto en común?



# ¿Cómo representas un texto con números?

**Asignando a cada carácter un valor**



**Asignando a cada palabra un valor o posición de un vector**



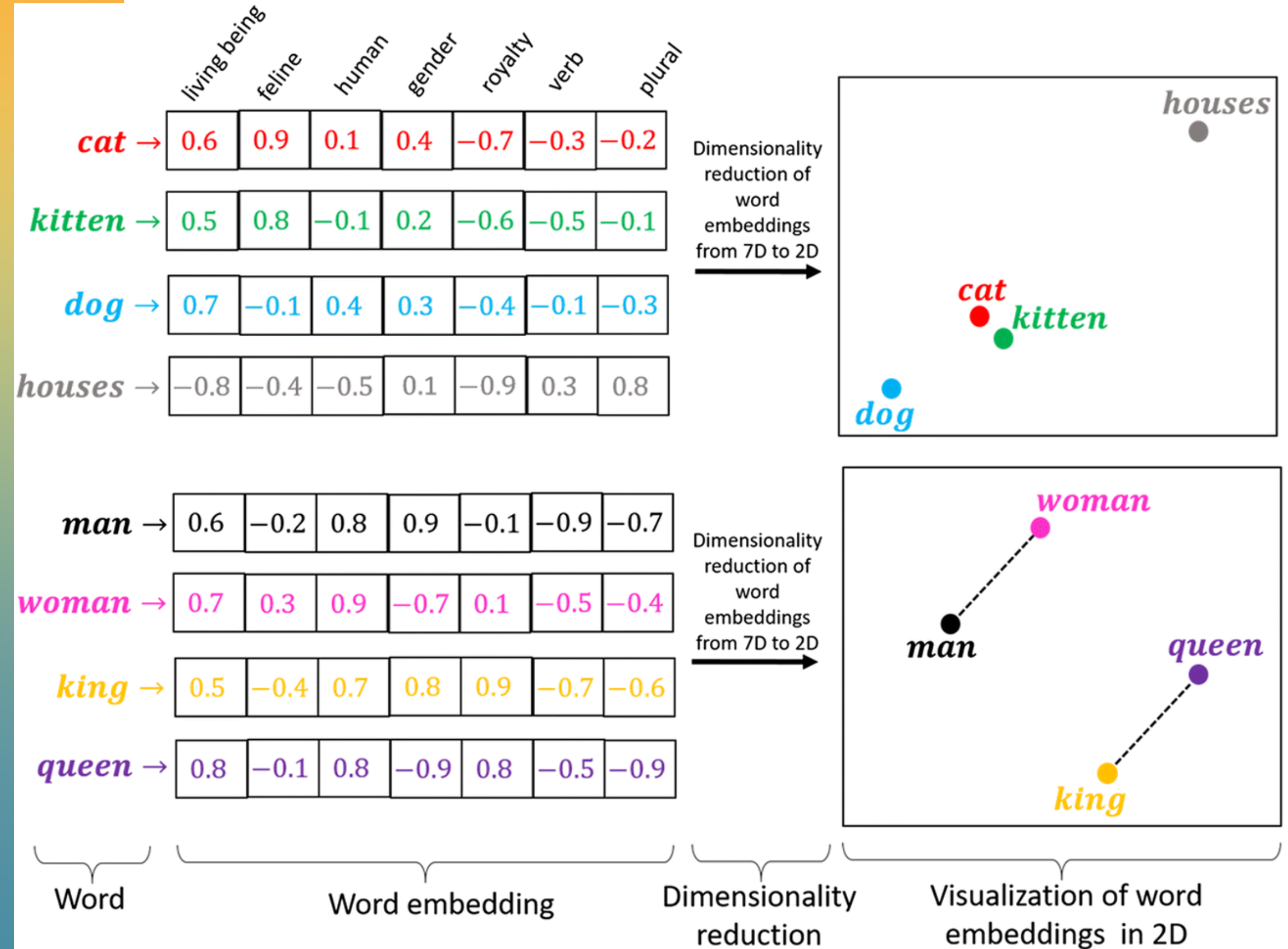
**Vectores para representar las palabras y sentido semántico**





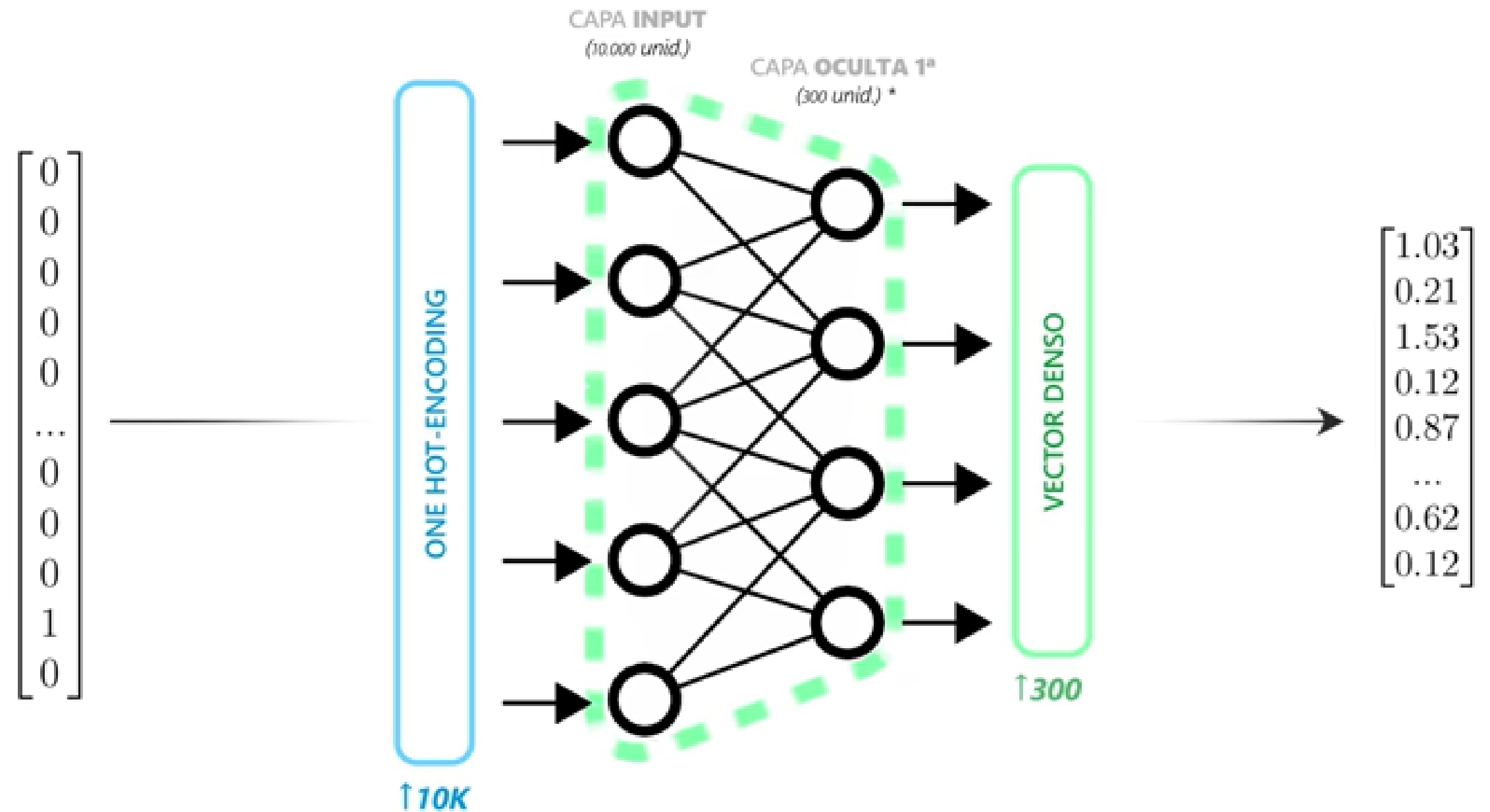
# Embeddings

Reducir la dimensionalidad



# Embeddings

Reducir la dimensionalidad



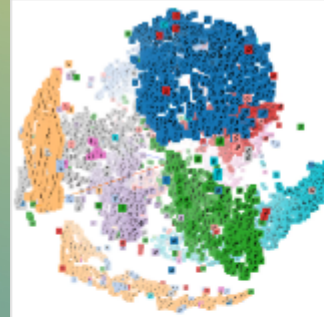
Tomado de: [https://www.youtube.com/watch?v=RkYuH\\_K7Fx4](https://www.youtube.com/watch?v=RkYuH_K7Fx4)

# Embeddings

Reducir la dimensionalidad

# Word2Vec

Google - 2013



**Embedding projector - visualization of high-dimensional data**

Visualize high dimensional data.

 [tensorflow.org /](https://www.tensorflow.org/)



# Preprocesamiento de texto



# Preprocesamiento

“Hey Amazon - my package never arrived  
[https://www.amazon.com/gp/css/order-history?ref\\_=nav\\_orders\\_first](https://www.amazon.com/gp/css/order-history?ref_=nav_orders_first)  
PLEASE FIX ASAP! @amazonhelp”

**Normalizar**

“hey amazon my package  
never arrived please fix asap”

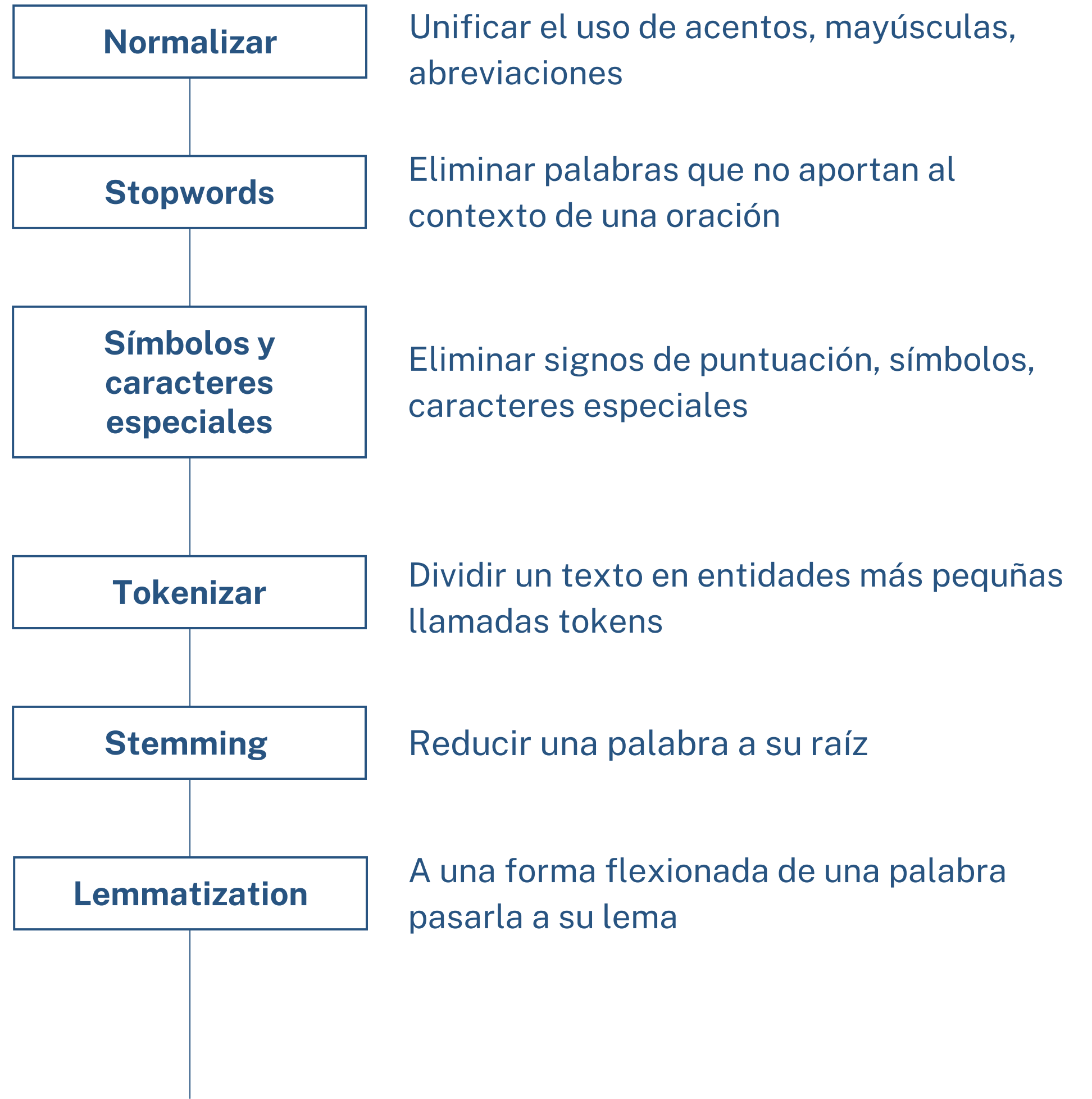
**Eliminar stopwords**

“amazon package never arrived fix  
asap”

**Tokenizar**

["amazon", "package",  
"never", "arrived", "fix",  
"asap"]

# Preprocesamiento

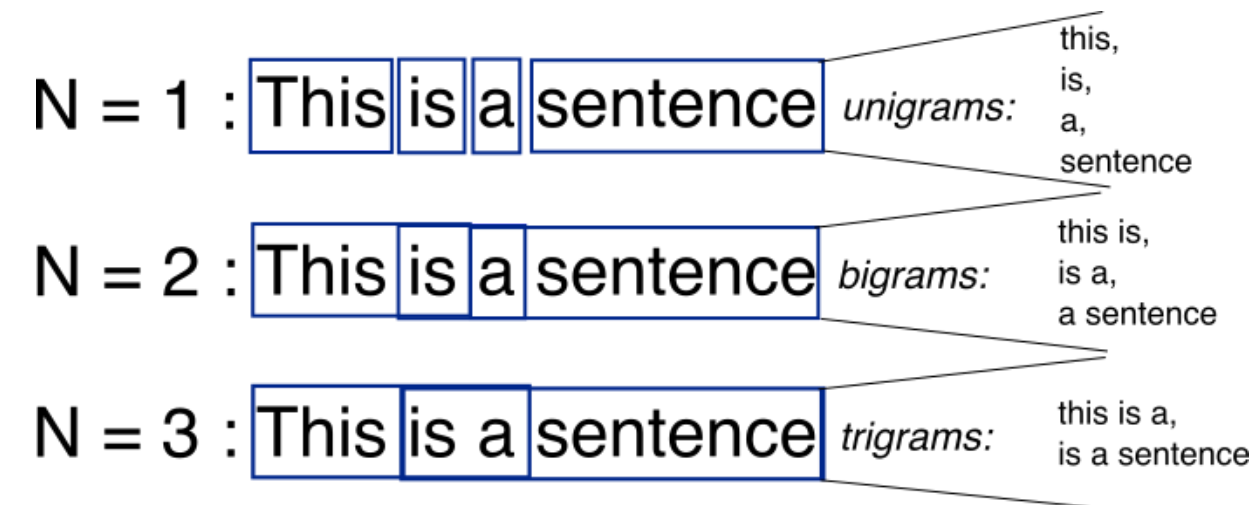


# Aplicaciones

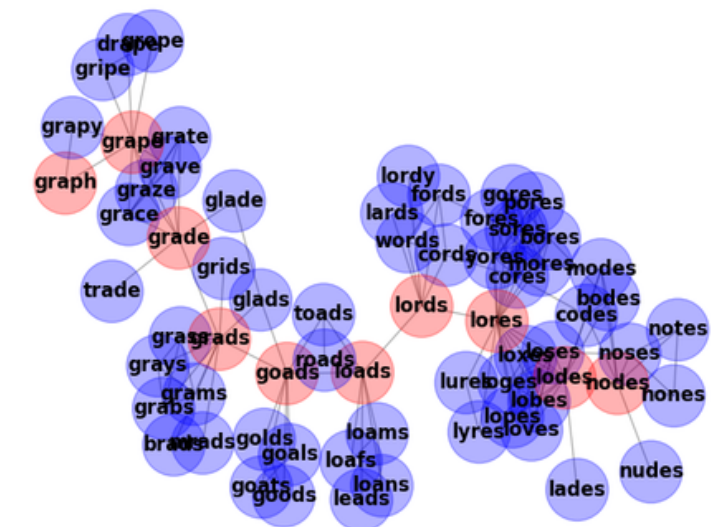
# ¿Qué podemos hacer con NLP?



## Wordcloud



## N-Grams



## Grafos de palabras

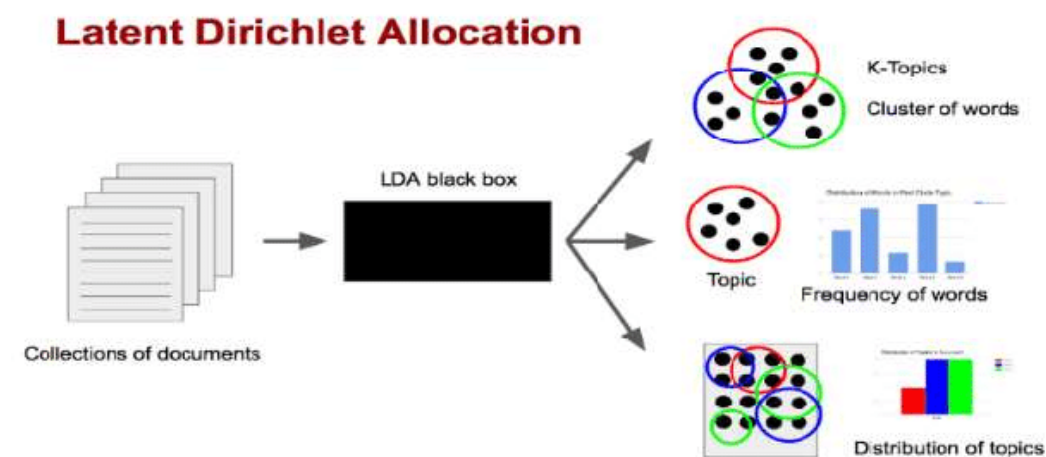


# ¿Qué podemos hacer con NLP?

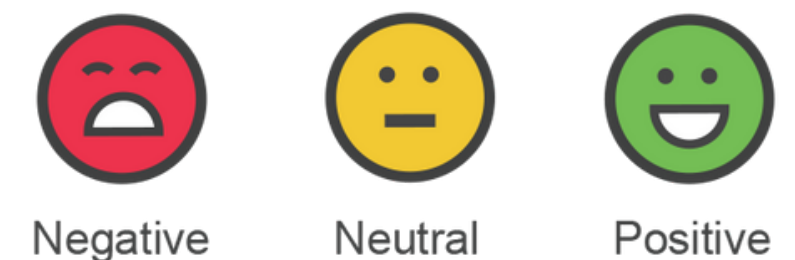
When **Sebastian Thrun** PERSON started at **Google** ORG in **2007** DATE, few people outside of the company took him seriously. "I can tell you very senior CEOs of major **American** NORP car companies would shake my hand and turn away because I wasn't worth talking to," said **Thrun** PERSON, now the co-founder and CEO of online higher education startup Udacity, in an interview with **Recode** ORG **earlier this week** DATE.

A little **less than a decade later** DATE, dozens of self-driving startups have cropped up while automakers around the world clamor, wallet in hand, to secure their place in the fast-moving world of fully automated transportation.

Named Entity Recognition (NER) y Part Of Speech Tagging (POS)



Latent Dirichlet Allocation (LDA)



Análisis de sentimientos y emociones

Vamos al código con  
algo de PySpark



# ¿Tienes alguna pregunta?

[github.com/lauralpezb/NLP\\_MyFutureAI](https://github.com/lauralpezb/NLP_MyFutureAI)

@lauralpezb