

- TF-idf \rightarrow Frecuencias

\rightarrow contador, lleva la cuenta en el corpus

\rightarrow Matriz resultante depende de $|V|$ donde $|V|$: Tamaño del vocabulario

\rightarrow "Maldición de la dimensionalidad"

poca data

muchas features (e.g. tfidf - 16000 columns)
"features"

Word Embeddings \rightarrow vector

"dog" \rightarrow

0.5	0.98	-0.1	0.25
-----	------	------	------

(dense) no ceros

$N=4$

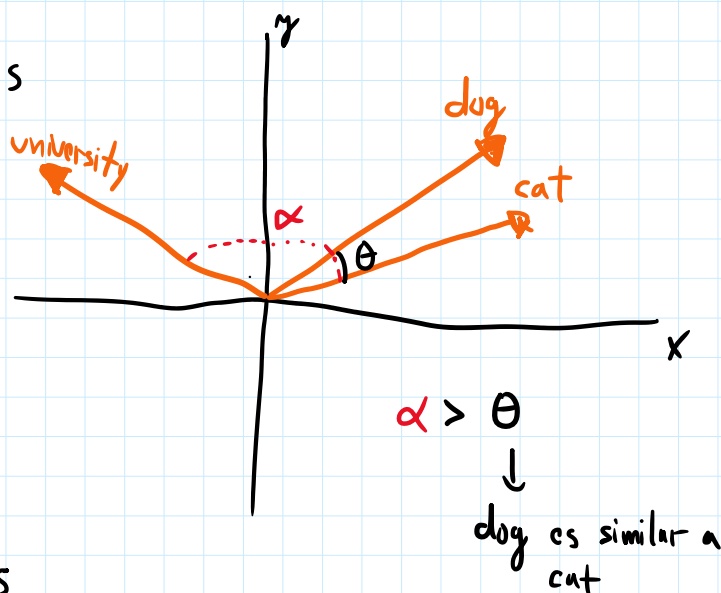
N es definida por nosotros

"cat" \rightarrow

0.1	-0.9	0.981	0.31
-----	------	-------	------

"university" \rightarrow

0.01	0.008	0.5	-0.21
------	-------	-----	-------



Word embeddings \rightarrow embedding estáticos

\rightarrow diccionario con repre. vectoriales de cada palabra en V .

modelo["cat"] =

--	--	--	--

modelo["dog"] =

--	--	--	--

\Leftrightarrow "cat" y "dog" $\in V$

\rightarrow dinámicos (Transformers (e.g. BERT))

→ dinámicos (Transformers (e.g. BERT))

↳ SOTA en NLP

Word2Vec → word → vector (estático)

- ↳ NO es un algoritmo/modelo
- ↳ ES una librería/paquete
- ↳ Skipgram, CBOW

Skipgram:

D1: [Hola, mundo, como, están, en, el], planeta, tierra

w_1 w_2 w_3 w_4 w_5 w_6 w_7 w_8

↑ ↑ ↑ ↑

$p(w_4, w_3) = p(+|w, c)$

$= p(-|w, c) = 1 - p(+|w, c)$

$p(+|w, c) = \text{similitud}(w, c)$

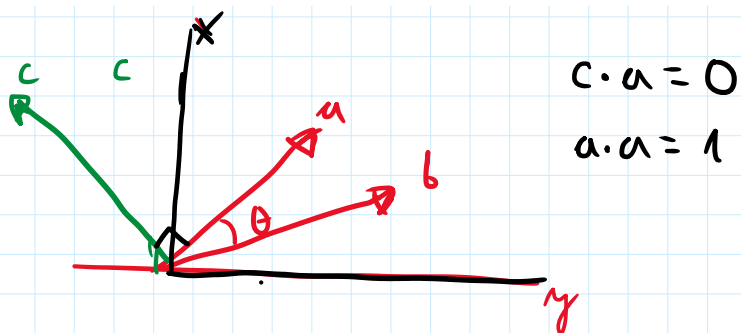
↳ producto punto (dot product) $\in \mathbb{R}$

$$\langle a, b \rangle \cdot \langle c, d \rangle = a \cdot c + b \cdot d = \mathbb{R}$$

$$\vec{a} \cdot \vec{b} = \cos \theta \rightarrow \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|} = \cos \theta \Rightarrow \vec{a} \cdot \vec{b} = |\vec{a}| |\vec{b}| \cos \theta$$

$\theta = \text{ángulo entre } \vec{a} \text{ y } \vec{b}$

$$|\vec{v}| = \sqrt{v_1^2 + v_2^2 + \dots + v_n^2} \in \mathbb{R}$$



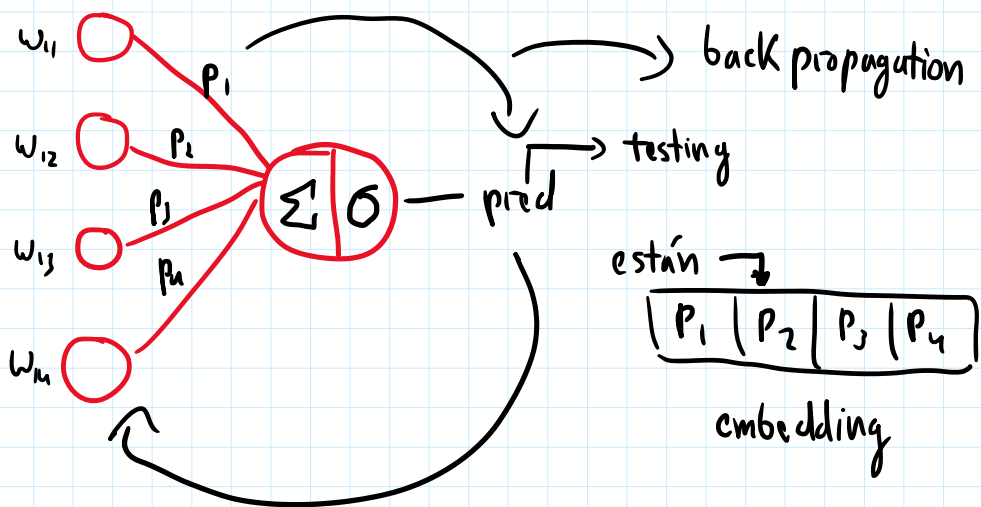
$$p(+|w, c) = \sigma(w \cdot c) \quad \sigma(x) = \frac{1}{1 + e^{-x}} \Rightarrow \text{sigmoide}$$

$$p(+|w, c) = \prod_{i=1}^L \sigma(w \cdot c_i) \quad \text{donde } L=4 \quad \text{window}=4$$

$$\sigma(w \cdot c_1) \cdot \sigma(w \cdot c_2) \cdot \dots$$

están \rightarrow

0.1	0.3	0.8	0.4
-----	-----	-----	-----



extraer muestras negativas \rightarrow Negative sampling

Skipgram with Negative sampling \rightarrow regresión logística (sigmoide)

OOV (Out-of-Vocabulary) \rightarrow Word2Vec (problema)

- \rightarrow random vector
- \rightarrow random (seed) \rightarrow fijas
- \rightarrow otro modelo \rightarrow FastText (n-grams)

"Hola mundo ESAN"

$n=2$ (bi-grams)

[Hola_mundo], [mundo-ESAN]

$n=2$

Hola

[<H, ol, a>] \rightarrow ola

\downarrow

ol + a = ola \rightarrow [vector]

$D_1 = [w_1, w_2, w_3] \rightarrow$

--	--	--	--

$w_1 \rightarrow$

--	--	--	--

$w_2 \rightarrow$

--	--	--	--

$w_3 \rightarrow$

--	--	--	--

$=$

\rightarrow

--	--	--	--

+

--	--	--	--

+

--	--	--	--

3

\rightarrow SUM (

--	--	--	--

)



SUM $\left(\begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} \begin{array}{l} D_1 \\ D_2 \\ D_3 \end{array} \right) \rightarrow \text{vector} \checkmark$

SUM $\left(\begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} \begin{array}{l} D_1 \\ w_1 \in D_1 \\ w_2 \in D_1 \end{array} \right) \checkmark$

$$\text{SVM} \left(\begin{array}{c|c} \text{ } & w_i \in D_i \\ \hline \text{ } & w_2 \in D_i \\ \hline \text{ } & w_3 \in D_i \\ \hline \end{array} \right) \quad \times$$

\hookrightarrow a vector