

Procesamiento de Lenguaje Natural

Olivia Gutú y Julio Weissman

Maestría en Ciencia de Datos
Semana 3: Modelos de espacios vectoriales



Luis estudia solo piano
Luis estudia solo una materia
↑
diferente significado

no me gusta practicar
odio hacer ejercicios
↑
mismo significado

Ventana de distancia k de un token: tokens que se encuentran a una distancia menor o igual a k .

los cronopios según Cortázar cronopios y famas de Cortázar

e. g.

Frecuencia w-w en una ventana de tamaño $k = 2$:

	los	cronopios	según	y	famas	de
Cortázar	0	2	1	1	1	1

Cortázar tendría un representación en $\mathbb{R}^{|V|}$, V el vocabulario de tipos.

Frecuencia de un token aparece en alguna categoría:

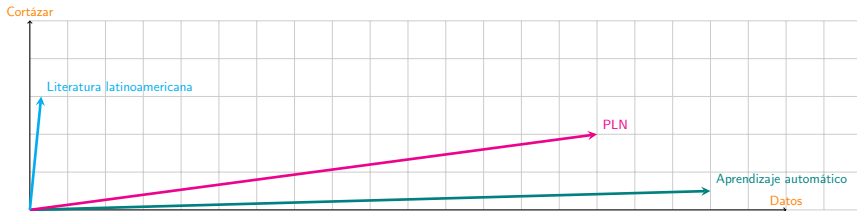
	Aprendizaje automático	PLN	Literatura latinoamericana
Datos	345445	233341	102
Cortázar	207	2345	34456

Podemos ver al vocabulario en términos de las categorías, o bien, las categorías en términos del vocabulario.

Distancia entre palabras



Universidad de Sonora



Distancia entre palabras: producto punto



Universidad de Sonora

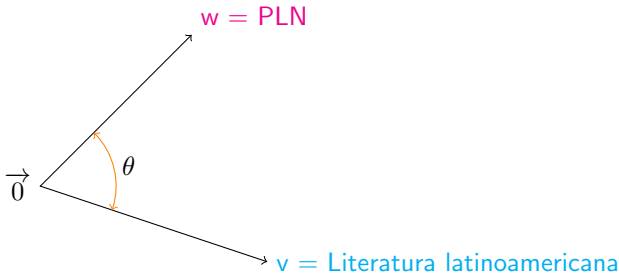
$$v = (102, 34456)$$

$$w = (233341, 2345)$$

$$v \text{ dot } w = (102 \times 233341) + (34456 \times 2345)$$

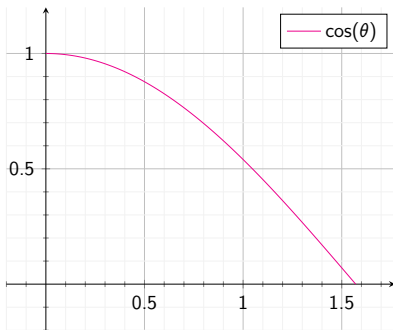
$$\begin{aligned} v \text{ dot } w &\leftarrow \text{producto punto de } v \text{ y } w \\ \|v\| = \sqrt{v \text{ dot } v} &\leftarrow \text{norma euclidiana de } v \\ \|v - w\| &\leftarrow \text{distancia euclidiana entre } v \text{ y } w \end{aligned}$$

$$\text{sim}(\mathbf{v}, \mathbf{w}) = \frac{\mathbf{v} \cdot \mathbf{w}}{\|\mathbf{v}\| \|\mathbf{w}\|}$$



¡justamente es el coseno del ángulo entre los dos vectores!

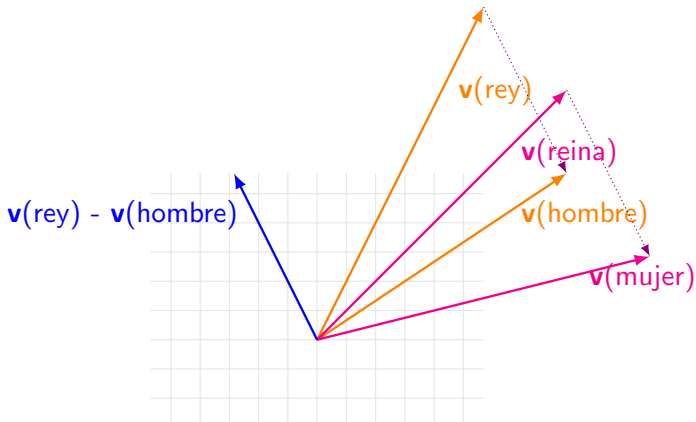
- La distancia euclidiana no es adecuada si los corpus son de tamaños diferentes.
- $\text{sim}(\mathbf{v}, \mathbf{w}) = 0$ implica que \mathbf{v} y \mathbf{w} son ortogonales: son lo menos parecido posible.
- $\text{sim}(\mathbf{v}, \mathbf{w}) = 1$ implica que \mathbf{v} y \mathbf{w} son linealmente dependientes: son lo más similar posible.



Distancia entre palabras: inferencia



Universidad de Sonora



$$v(\text{rey}) - v(\text{hombre}) + v(\text{mujer}) \rightarrow v(\text{reina})$$

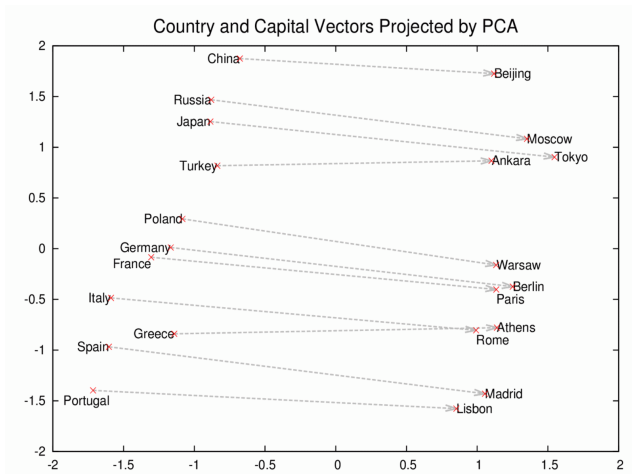


Figura tomada del artículo original de Mikolov et. al (2014)

