



Procesamiento de Lenguaje Natural

Vectores Semánticos

Mauricio Toledo-Acosta
mauricio.toledo@unison.mx

Departamento de Matemáticas
Universidad de Sonora



Section 1

Introducción



Procesamiento de Lenguaje Natural

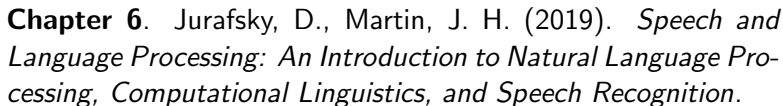
Introducción

Similitud
coseno

Bag of Words

Tf-Idf

The screenshot shows a presentation slide with a dark background. At the top left is the 'the.language.nerds' logo, which consists of a circular icon with a brain and the text 'the.language.nerds'. To the right of the logo is the text 'the.language.nerds' in a white sans-serif font. In the top right corner, there are three vertical white dots. The main content of the slide is the title 'Me speaking English:' in a large, bold, yellow serif font. Below the title is a paragraph of text in a white serif font: 'I have literally no idea what this word translates to in my native language but I've seen it being used in similar context so I'm just gonna use it here and pray that it does mean what I think it means.'





Palabras que aparecen en contextos similares tienden a tener significados similares. Este vínculo entre la similitud en la distribución de las palabras y la similitud en su significado se denomina **hipótesis distribucional**.

Esta hipótesis fue formulada en los años 50 por lingüistas como Joos (1950), Harris (1954) y Firth (1957), que observaron que las palabras sinónimas (como oculista y oftalmólogo) tendían a aparecer en el mismo entorno (por ejemplo, cerca de palabras como ojo o examinado).

A word is characterized by the company it keeps

Más información



Ejemplo

Procesamiento
de Lenguaje
Natural

Introducción

Similitud
coseno

Bag of Words

Tf-Idf

¿Qué es el **Ongchoi**?



- **Ongchoi** es delicioso **salteado con ajo**.
- **Ongchoi** es excelente **sobre arroz**.
- ...**ongchoi** **hojas** con salsas saladas...

- ...la espinaca salteada con ajo sobre arroz...
- ...los tallos de acelga y las hojas son deliciosos...
- ...la col rizada y otras verduras de hoja saladas...

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻ 6/29



La cercanía de embeddings da cuenta de diversos fenómenos, además de la similitud de palabras.

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻ 7/29



◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻ 8/29



Similitud de Palabras

Procesamiento
de Lenguaje
Natural

Introducción

Similitud
coseno

Bag of Words

Tf-Idf

No todas las palabras tienen muchos sinónimos, sin embargo, la mayoría de ellas tienen muchas palabras similares (gato – perro). La noción de similitud entre palabras es muy útil en diversas tareas semánticas, por ejemplo, decidir si dos oraciones significan cosas parecidas.

- Podemos obtener las similitudes entre palabras de listas pre-definidas (por ejemplo, [SimLex-999](#)).

vanish	disappear	9.8
behave	obey	7.3
belief	impression	5.95
muscle	bone	3.65
modest	flexible	0.98
hole	agreement	0.3

- Podemos aprender las similitudes a partir de co-ocurrencias.



Modelos de semántica vectorial

Procesamiento
de Lenguaje
Natural

Introducción

Similitud
coseno

Bag of Words

Tf-Idf

- Modelos basados en conteos: El significado de una palabra está dado en términos de ocurrencias en documentos.
 - Bag of Words (BOW)
 - Term Frequency - Inverse Document Frequency (TF-IDF)
- Modelos basados en redes neuronales:
 - Clásicos: Word2Vec, GloVe, ...
 - LLMs: GPT, LLaMA, DeepSeek, Gemma, Qwen,...



- Modelos basados en conteos:
 - + Modelo sencillo y simple
 - + Interpretabilidad
 - Vectores raros (*sparse*)
 - Alta dimensionalidad (del orden de miles o más)
- Modelos basados en redes neuronales:
 - + Vectores densos
 - + Menor dimensionalidad (del orden de cientos)
 - Algunos pueden ser computacionalmente caros de obtener



¿Cómo medimos la similitud/distancia entre vectores?

Consideremos dos embeddings $u, v \in \mathbb{R}^D$.

- La **similitud coseno** es un valor entre -1 y 1 y está dada por

$$\text{sim}(u, v) = \frac{u \cdot v}{\|u\| \|v\|} = \cos(\angle u, v)$$

donde $u \cdot v$ es el producto punto y $\|u\|$ es la norma del vector. Podemos normalizar $\|u\| = \|v\| = 1$ y tenemos

$$\text{sim}(u, v) = u \cdot v$$

En general, no usamos la distancia Euclidiana.

Ejemplo



Consideremos dos embeddings $u, v \in \mathbb{R}^D$.

- La distancia angular es un valor entre 0 y π dado por

$$d_{\theta}(u, v) = \arccos(\text{sim}(u, v))$$

- En ocasiones, nos referimos a la métrica coseno como

$$d_{\cos}(u, v) = 1 - \text{sim}(u, v)$$

En general, no usamos la distancia Euclidiana.

Ejemplo



La matriz term-document

Procesamiento
de Lenguaje
Natural

Introducción

Similitud
coseno

Bag of Words

Tf-Idf

- La Revolución Francesa fue un período de grandes **cambios** políticos y sociales en **Europa**.
- El Imperio Romano dominó gran parte de **Europa** durante siglos, expandiéndose por toda **Europa**.
- La paella es un plato tradicional de España que lleva **arroz**, mariscos y verduras.
- El sushi es una comida japonesa hecha con **arroz** y **pescado** crudo, acompañado de algas.

Texto	Europa	cambios	arroz	pescado
1	1	1	0	0
2	2	0	0	0
3	0	0	1	0
4	0	0	1	1

Visualización



Modelo BOW

Procesamiento
de Lenguaje
Natural

Introducción

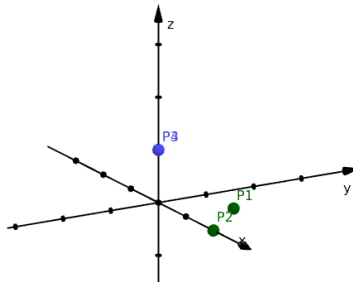
Similitud
coseno

Bag of Words

Tf-Idf

Texto	Europa	cambios	arroz	pescado
1	1	1	0	0
2	2	0	0	0
3	0	0	1	0
4	0	0	1	1

El modelo BOW asigna a cada documento el vector correspondiente a la fila. El vector de cada palabra es su columna.





Preguntas

Procesamiento
de Lenguaje
Natural

Introducción

Similitud
coseno

Bag of Words

Tf-Idf

- ¿Qué obtenemos si sumamos las filas?
- ¿Qué obtenemos si sumamos las columnas?
- ¿Qué tamaño tiene la matriz anterior?
- ¿Qué palabras tenderán a dominar la matriz?
- ¿Qué significa que dos palabras sean similares usando estas representaciones vectoriales?



¿Qué significa que dos palabras sean similares usando estas representaciones vectoriales?

Las palabras similares tienen vectores similares porque suelen aparecer en documentos similares. La matriz término-documento nos permite representar el significado de una palabra por los documentos en los que suele aparecer.



la matriz término-término (también denominada matriz palabra-palabra o matriz término-contexto) es la matriz de tamaño $|V| \times |V|$ donde las columnas y filas están etiquetadas por palabras. La entrada ij es el número de veces que la palabra i (objetivo) y la palabra j (contexto) coinciden en algún contexto (ventana) en algún corpus de entrenamiento.

Doc₁: *I go to school every day by bus.*

Doc₂: *I go to theatre every night by bus.*



Las matrices término-documento se definieron originalmente como un medio de encontrar documentos similares para la tarea de recuperación de información documental. Dos documentos que son similares tenderán a tener palabras similares, y si dos documentos tienen palabras similares sus vectores columna tenderán a ser similares.

IR es la tarea que consiste en buscar, localizar y presentar información que coincida con la consulta de búsqueda o la necesidad de información de un usuario.



Doc₃: *Las playas en México tienen agua templada...*

	Doc ₁	Doc ₂	Doc ₃
agua	1	2	1
ingredientes	1	2	0
...



Doc₃: *Las playas en México tienen agua templada...*

	Doc ₁	Doc ₂	Doc ₃
agua	1	2	1
ingredientes	1	2	0
...

Con la métrica euclidiana, el más similar a Doc_1 es Doc_3 . Con la métrica coseno, es Doc_2 .



Documentos:

Palabras consideradas:



Section 4

Tf-Idf





Doc₃: *Las playas en México tienen agua templada...*

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻ 27/29



Ejemplo

Procesamiento
de Lenguaje
Natural

Introducción

Similitud
coseno

Bag of Words

Tf-Idf

BOW:

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

TF-IDF:

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	0.074	0	0.22	0.28
good	0	0	0	0
fool	0.019	0.021	0.0036	0.0083
wit	0.049	0.044	0.018	0.022



Aplicaciones Adicionales

Procesamiento de Lenguaje Natural

Introducción

Similitud coseno

Bag of Words

Tf-Idf

Una aplicación adicional de estos modelos es la extracción de features del texto para tareas de Machine Learning. Es importante reflexionar sobre qué rasgos del texto captan estas features.

- Clasificación (Análisis de Sentimientos, etc.).
- Segmentación (Detección de tópicos, etc.).
- Similitud de documentos (IR, Autoría, etc.).