



# Procesamiento de Lenguaje Natural

## Vectores Semánticos

Mauricio Toledo-Acosta  
[mauricio.toledo@unison.mx](mailto:mauricio.toledo@unison.mx)

Departamento de Matemáticas  
Universidad de Sonora



Procesamiento  
de Lenguaje  
Natural

Introducción

Modelo BOW

Modelo  
TF-IDF

La métrica  
angular

## Section 1

### Introducción



# Vectores Semánticos

Representación numérica del lenguaje natural

Procesamiento  
de Lenguaje  
Natural

Introducción

Modelo BOW

Modelo  
TF-IDF

La métrica  
angular

## ¿Por qué representar texto numéricamente?

- Los algoritmos de ML requieren features numéricas
- Necesidad de cuantificar la semántica y similitud
- Búsqueda y recuperación eficiente de información
- Análisis automatizado de grandes volúmenes de texto



# Vectores Semánticos

Representación numérica del lenguaje natural

## Evolución de las representaciones

- **Modelos tradicionales:** BOW, TF-IDF (representaciones locales)
  - **Modelos modernos:** Word2Vec, GloVe, BERT (representaciones densas)
  - **LLMs:** GPT, LLaMa, Qwen, Claude, DeepSeek.
- 
- De representaciones dispersas a embeddings densos
  - De contar palabras a capturar significado contextual



## Section 2

### Modelo BOW



# La matriz term-document

Procesamiento  
de Lenguaje  
Natural

Introducción

Modelo BOW

Modelo  
TF-IDF

La métrica  
angular

- 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



# Modelo BOW

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.

¿Cuántas filas y cuántas columnas hay en una matriz BOW?

Visualización



# Ventajas y Desventajas

## Ventajas:

- Simplicidad e interpretabilidad
- Fácil implementación computacional
- Eficiente para conjuntos de datos grandes
- Base fundamental para modelos más avanzados
- Compatible con algoritmos de machine learning

## Desventajas:

- Pérdida del orden de las palabras, contexto y semántica
- No captura relaciones entre términos (sinonimia, polisemia, etc.)
- Matriz grande y *sparse*
- Sensible a stop words
- No considera la importancia relativa de términos



# Ejercicio BOW

Procesamiento  
de Lenguaje  
Natural

Introducción

Modelo BOW

Modelo  
TF-IDF

La métrica  
angular

## Documentos:

- El gato come ratones y juega con el perro. El perro duerme al lado y come.
- El gato come pescado.
- El perro ladra fuerte y come.
- El código tiene un error.
- El programa ejecuta código.

## Palabras consideradas:

- |           |           |            |
|-----------|-----------|------------|
| ● gato    | ● duerme  | ● código   |
| ● come    | ● lado    | ● error    |
| ● ratones | ● pescado | ● programa |
| ● juega   | ● ladra   | ● ejecuta  |
| ● perro   | ● fuerte  |            |



## Section 3

### Modelo TF-IDF



# El modelo TF-IDF

Problemas de la representación BOW (por conteos):

- Las palabras muy comunes ("el", "de", "y") dominan los vectores. No necesariamente las stopwords.
- Estas palabras aportan poco sobre el contenido real del documento.

Idea central del modelo TF-IDF:

- No todas las palabras son igual de importantes.
- Una palabra es relevante si aparece *mucho* en un documento pero aparece en *pocos* documentos del corpus.



# La matriz 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	0.043	0.301	0	0
2	0.043	0	0	0
3	0	0	0.090	0
4	0	0	0.060	0.301

Los valores TF-IDF ponderan la importancia de cada término según su frecuencia en el documento y su rareza en el corpus.



# Cálculo del TF-IDF

- $\text{TF-IDF} = \text{TF} * \text{IDF}$
- TF (Term Frequency): Frecuencia del término en el documento
- IDF (Inverse Document Frequency): Rareza del término en el corpus

$$TF(t, d) = \frac{\text{frecuencia del término } t \text{ en documento } d}{\text{total de términos en documento } d}$$

$$IDF(t) = \log \left( \frac{\text{total de documentos}}{\text{documentos que contienen término } t} \right)$$

$$TF-IDF(t, d) = TF(t, d) * IDF(t)$$



# Cálculo del TF-IDF

## Ejemplo: *Europa* en Texto 2

- $TF = 2/14 = 0.143$  (aparece 2 veces de 14 palabras totales)
- $IDF = \log(4/2) = \log(2) = 0.301$
- $TF-IDF = 0.143 * 0.301 = 0.043$

## Ejemplo: *arroz* en Texto 3

- $TF = 1/9 = 0.111$  (aparece 1 vez de 9 palabras totales)
- $IDF = \log(4/1) = \log(4) = 0.602$
- $TF-IDF = 0.111 * 0.602 = 0.067$



# Ventajas y Desventajas

## Ventajas:

- Mejora la representación al ponderar términos
- Reduce el impacto de términos muy frecuentes
- Identifica términos relevantes por documento
- Mejor rendimiento en recuperación de información
- Fácil de implementar y computacionalmente eficiente

## Desventajas:

- No captura relaciones semánticas entre términos
- Matriz grande y *sparse*
- Sensible a la longitud del documento
- No resuelve problemas de polisemia o sinonimia
- Depende de la calidad del corpus de entrenamiento



## Section 4

### La métrica angular



# ¿Cómo medir la distancia entre documentos?

- Doc 0: “buen hotel, economico, muy recomendable hotel, mi hotel favorito”
- Doc 1: “buen hotel, economico”
- Doc 2: “buen automovil economico”

Representamos cada documento como vector BOW.

	automovil	buen	economico	favorito	hotel	recomendable
0	0	1	1	1	3	1
1	0	1	1	0	1	0
2	1	1	1	0	0	0

¿A qué documento se parece más el Doc 1? Calculemos distancias



# ¿Cómo medir la distancia entre documentos?

- Con distancia Euclidiana: Doc 1 es más cercano a Doc 2.

La distancia euclidiana es sensible a la longitud de los documentos: Doc 0 es más largo.

- Con la distancia angular (coseno):

$$d_\theta(u, v) = \arccos \left( \frac{u \cdot v}{|u| \cdot |v|} \right)$$

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

Doc 1 es más cercano a Doc 0.



# Conclusión

Procesamiento  
de Lenguaje  
Natural

Introducción

Modelo BOW

Modelo  
TF-IDF

La métrica  
angular

La similitud coseno, o métrica angular, suele ser la métrica ideal para medir distancias entre representaciones de documentos (incluso en modelos basados en redes neuronales).

Dos documentos son similares si sus vectores apuntan en una dirección similar, sin importar su longitud. La distancia coseno se centra en el contenido de palabras. De esta forma, ignora la longitud del documento y captura mejor la similitud temática.