



# Procesamiento de Lenguaje Natural

## Vectores Semánticos

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

Departamento de Matemáticas  
Universidad de Sonora



## Section 1

# Introducción



## Section 2

# Ejemplos ilustrativos



# La matriz term-document

Procesamiento  
de Lenguaje  
Natural

Introducción

Ejemplos  
ilustrativos

Modelo TF-IDF

LSA

- 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

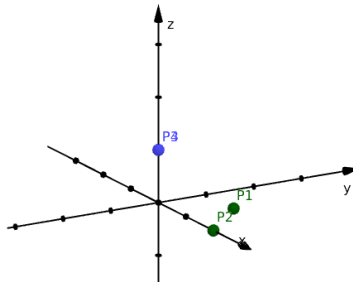
Ejemplos  
ilustrativos

Modelo TF-IDF

LSA

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.





- 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.

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



## Subsection 1

# Modelo TF-IDF



# 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	0.301	0.602	0	0
2	0.602	0	0	0
3	0	0	0.301	0
4	0	0	0.301	0.602

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

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

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

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

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

Ejemplo de *Europa* en Texto 2:

- $TF = 2/7 = 0.286$  (aparece 2 veces de 7 palabras totales)
- $IDF = \log(4/2) = \log(2) = 0.301$
- $TF-IDF = 0.286 * 0.301 = 0.086$



## Section 3

# LSA



# Latent Semantic Analysis

Procesamiento  
de Lenguaje  
Natural

Introducción

Ejemplos  
ilustrativos

Modelo TF-IDF

LSA

## LSA (Latent Semantic Analysis)

Técnica de procesamiento de lenguaje natural usada en Topic Modelling para descubrir temas en textos, es decir, identificar temas ocultos en un conjunto de documentos.



- **Matriz Término-Documento:** Representación numérica de textos, típicamente BOW o TF-IDF.
- **SVD:** Reducción de dimensionalidad para capturar relaciones semánticas.
- **Espacio semántico latente:** Representación compacta de palabras y documentos.



- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻ 13/19



# Espacio semántico latente

Procesamiento  
de Lenguaje  
Natural

Introducción

Ejemplos  
ilustrativos

Modelo TF-IDF

LSA

- **Concepto:** Representación de palabras y documentos en un espacio de menor dimensión.
- **Ventaja:** Captura relaciones semánticas entre términos y documentos.
- **Ejemplo:** Palabras como "coche" y "automóvil" estarán cerca.



# Proceso de LSA

## Procesamiento de Lenguaje Natural

### Introducción

### Ejemplos ilustrativos

Modelo TF-IDF

### LSA

- **Preprocesamiento:** Tokenización, eliminación de stopwords, etc.
- **Matriz Término-Documento:** Creación y ponderación (TF-IDF).
- **SVD:** Aplicación y reducción de dimensionalidad.
- **Interpretación:** Identificación de temas latentes.



# Ventajas de LSA

Procesamiento  
de Lenguaje  
Natural

Introducción

Ejemplos  
ilustrativos

Modelo TF-IDF

LSA

- Captura relaciones semánticas entre palabras.
- Reduce el ruido en grandes conjuntos de datos.
- Simple y fácil de implementar.





# Limitaciones de LSA

## Procesamiento de Lenguaje Natural

### Introducción

### Ejemplos ilustrativos

Modelo TF-IDF

### LSA

- Dificultad para interpretar temas explícitamente.
- Depende del preprocesamiento y parámetros.



# Aplicaciones de LSA

Procesamiento  
de Lenguaje  
Natural

Introducción

Ejemplos  
ilustrativos

Modelo TF-IDF

LSA

- Recuperación de información.
- Clasificación de textos.
- Análisis de sentimientos.
- Recomendación de contenido.



- 19/19