



UNIVERSIDAD CENTROCCIDENTA
"LISANDRO ALVARADO"
DECANATO DE CIENCIAS Y TECNOLOGÍA



Laboratorio unidad IV

Introducción a la Inteligencia Artificial



Sección 1

V-27.411.922, Castillo Gustavo

V-25.401.656, Colmenárez Luis

V-28.454.911, Rojas Gabriel

Profesora: María Perez

Barquisimeto, diciembre de 2024

Índice

Resumen	3
Introducción	4
Objetivos de la investigación	5
Marco Teórico	6
1. Aprendizaje profundo (Deep Learning)	6
2. Modelo de memoria distribuida (Paragraph Vector)	6
3. ¿Qué es dataset?	6
3.1 Descripción de los conjuntos de datos del trabajo:	7
4. ¿Qué es Gensim y sus usos?	7
5. Modelo preentrenado Doc2Vec	7
6. Diferencias entre Doc2Vec y Word2Vec	8
7. ¿Cómo entrenar nuestros propios vectores de palabras usando Gensim Doc2Vec?	8
Resultados	9
1. Experimentar con wikipedia.txt	9
2. Experimentar con boletines.csv	13
3. Experimentar con boletinesprueba.csv	18
Conclusiones	19
Referencias Bibliográficas	21
Referencias Web	

Resumen

Este informe documenta el entrenamiento y evaluación de un modelo de aprendizaje automático basado en la técnica Doc2Vec de la biblioteca Gensim, utilizando tres conjuntos de datos: Wikipedia, boletines, y boletines de prueba. La metodología consistió en procesar los datos textuales mediante técnicas de limpieza como la eliminación de puntuación y palabras irrelevantes, para luego etiquetar y construir representaciones vectoriales que preserven el significado semántico de los documentos.

El modelo se entrenó en iteraciones ajustando parámetros como la tasa de aprendizaje. Con los datos de Wikipedia, se logró construir un modelo genérico para representar texto general. Posteriormente, los conjuntos de boletines y boletines de prueba fueron empleados para personalizar el modelo hacia el análisis de datos especializados, evaluando similitudes entre documentos.

Los resultados mostraron que el modelo puede identificar relaciones entre documentos con base en similitudes semánticas, lo que sugiere una representación adecuada de los datos. Las conclusiones destacan la eficacia del enfoque Doc2Vec para tareas de clasificación y agrupamiento de texto en dominios específicos, así como la importancia de datos limpios y variados en el entrenamiento para mejorar la generalización del modelo.

Introducción

En la era de la información, el análisis eficiente de datos textuales se ha convertido en una herramienta fundamental para diversas aplicaciones, como la búsqueda de información, la clasificación de documentos, y el análisis de contenido. Una de las técnicas más avanzadas en este campo es el modelo Doc2Vec, que permite representar documentos como vectores numéricos conservando relaciones semánticas.

Este informe detalla el proceso de entrenamiento y evaluación de un modelo Doc2Vec utilizando tres conjuntos de datos: definiciones generales extraídas de Wikipedia y textos especializados de boletines. El objetivo principal es explorar la capacidad del modelo para adaptarse a diferentes dominios y evaluar su desempeño al identificar relaciones semánticas en documentos previamente etiquetados.

Se describe la metodología utilizada para limpiar, procesar y etiquetar los datos, así como las configuraciones empleadas durante el entrenamiento. Finalmente, se presentan los resultados obtenidos en términos de similitudes documentales y las implicaciones prácticas del modelo para aplicaciones futuras en análisis de texto especializado.

Objetivos de la investigación

General:

Desarrollar y evaluar un modelo Doc2Vec para la representación semántica de documentos textuales en dominios generales y especializados.

Específicos:

- Implementar un modelo Doc2Vec utilizando datos generales de Wikipedia para explorar su capacidad de representar texto en contextos amplios.
- Entrenar y personalizar el modelo con datos específicos de boletines, evaluando su desempeño en tareas de agrupamiento y clasificación.
- Analizar la capacidad del modelo para identificar similitudes semánticas entre documentos en conjuntos de datos especializados y de prueba.
- Documentar la metodología y resultados obtenidos, resaltando las ventajas y limitaciones del enfoque Doc2Vec en el análisis de datos textuales.

Marco Teórico

1. Aprendizaje profundo (Deep Learning)

El aprendizaje profundo es una rama del aprendizaje automático que utiliza redes neuronales artificiales con múltiples capas para procesar grandes volúmenes de datos. Este enfoque busca simular el funcionamiento del cerebro humano, permitiendo que las máquinas realicen tareas como reconocimiento de imágenes, procesamiento del lenguaje natural y predicción basada en patrones complejos.

El aprendizaje profundo se destaca por su capacidad para automatizar la extracción de características relevantes de los datos, eliminando en muchos casos la necesidad de ingeniería manual de características. Esto lo convierte en una herramienta poderosa para análisis de texto, entre otros campos.

2. Modelo de memoria distribuida (Paragraph Vector)

El modelo de memoria distribuida, también conocido como Paragraph Vector, es una extensión del modelo Word2Vec diseñado para representar textos completos (frases, párrafos o documentos) como vectores de densidad fija. Este modelo asigna una representación vectorial única a cada documento que captura su contexto semántico global, permitiendo tareas como clasificación, búsqueda y agrupamiento.

A diferencia de los enfoques basados únicamente en palabras, Paragraph Vector considera la relación entre palabras y documentos en su conjunto, lo que lo hace ideal para analizar conjuntos de datos textuales heterogéneos.

3. ¿Qué es dataset?

Un dataset es un conjunto estructurado de datos que se utiliza para entrenar, validar y probar modelos de aprendizaje automático. Los datasets suelen estar organizados en

registros y atributos, donde cada registro representa una instancia única y los atributos describen sus características.

3.1 Descripción de los conjuntos de datos del trabajo:

- **Wikipedia.txt:** Contiene definiciones generales separadas por saltos de línea. Este dataset es utilizado para entrenar el modelo en un dominio amplio y no especializado.
- **Boletines.csv:** Incluye textos especializados en el formato Idtexto;Texto;Supervisión. Sirve como fuente de entrenamiento para adaptar el modelo a un dominio específico.
- **Boletines Prueba.csv:** Posee la misma estructura que Boletines.csv, pero se utiliza para evaluar el desempeño del modelo en un conjunto independiente.

4. ¿Qué es Gensim y sus usos?

Gensim es una biblioteca de Python diseñada para el procesamiento y análisis de texto, enfocándose en la modelización de temas, similitud semántica y representación de texto mediante vectores. Entre sus usos destacan:

- Entrenamiento de modelos como Word2Vec, Doc2Vec y FastText.
- Generación de representaciones vectoriales de texto para clasificación y agrupamiento.
- Similitud semántica y recuperación de información.

5. Modelo pre entrenado Doc2Vec

¿Qué es y cómo funciona?

Doc2Vec es una extensión de Word2Vec que permite representar documentos enteros como vectores en un espacio semántico. Su principio se basa en asignar un vector único a cada documento, el cual es entrenado simultáneamente con los vectores de palabras para capturar el contexto global del texto.

Principios:

Usa un vector adicional para representar el documento junto con los vectores de palabras.

Durante el entrenamiento, se ajustan los pesos del vector del documento para maximizar la probabilidad de predecir palabras dentro del texto.

Predicción del modelo:

Una vez entrenado, Doc2Vec utiliza el vector del documento para encontrar textos similares, agrupar documentos o clasificar textos.

6. Diferencias entre Doc2Vec y Word2Vec

Aspecto	Word2Vec	Doc2Vec
Representación	Palabras individuales.	Documentos completos.
Aplicaciones	Análisis léxico.	Clasificación, agrupamiento y similitud de documentos.
Contexto semántico	Basado en palabras vecinas.	Considera el documento completo.
Vectores únicos	No.	Sí.

7. ¿Cómo entrenar nuestros propios vectores de palabras usando Gensim Doc2Vec?

7.1. Preparación de los datos:

- Limpiar el texto eliminando puntuación, caracteres especiales y palabras irrelevantes.
- Dividir el texto en frases o documentos etiquetados.

7.2. Etiquetado de documentos:

- Usar estructuras como TaggedDocument de Gensim para asignar identificadores únicos a cada documento.

7.3. Construcción del modelo:

- Inicializar un modelo Doc2Vec configurando parámetros como la dimensión del vector y la tasa de aprendizaje.

7.4. Entrenamiento:

- Utilizar el método train en múltiples iteraciones, ajustando la tasa de aprendizaje para evitar sobreajuste.

7.5. Evaluación:

- Utilizar el método most_similar para verificar la capacidad del modelo de identificar documentos relacionados.

7.6. Guardado:

- Guardar el modelo entrenado con save para reutilización y análisis posterior.

Resultados

1. Experimentación y entrenamiento del modelo gensim_doc2vec con data Wikipedia – Pasos realizados, entrenamiento con la data de Wikipedia, printscreen de corrida y salidas...).

1.1. Pasos realizados

1.1.1. Carga de datos: La experimentación se realizó con un archivo de texto grande (wikipedia.txt) que contiene información en formato conceptual. Dado que el tamaño del archivo (561 MB) puede dificultar su procesamiento, se adoptaron estrategias para manejar

```

8 # Leer archivo y tomar una muestra aleatoria de líneas
9 file_path = "wikipedia.txt"
10 sample_size = 5000
11 with open(file_path, "r", encoding="UTF-8") as file:
12     lines = file.readlines() # Leer todas las líneas
13     sampled_lines = random.sample(lines, min(len(lines), sample_size)) # Seleccionar muestra aleatoria

```

el corpus de manera eficiente, utilizando una muestra representativa del archivo.

1.1.2. Preprocesamiento de datos:

- I. Se divide el texto en oraciones: Usando una expresión regular que identifica los finales de las oraciones.
- II. Divide las oraciones en palabras y se ignoran oraciones demasiado cortas.
- III. Se crea una etiqueta única para cada oración.
- IV. Se Almacena cada oración procesada en una lista de documentos.

```
15 # Preprocesar las líneas seleccionadas
16 sentenceEnders = re.compile(r'[.!?]\s+')
17 exclude = set(string.punctuation)
18 LabelDoc = namedtuple('LabelDoc', 'words tags')
19 all_docs = []
20 for idx, line in enumerate(sampled_lines):
21     sentences = sentenceEnders.split(line) # Dividir por signos de puntuación
22     for sentence in sentences:
23         sentence = ''.join(ch for ch in sentence if ch not in exclude) # Quitar puntuación
24         words = sentence.split() # Dividir en palabras
25         if len(words) < 3: # Ignorar frases muy cortas
26             continue
27         tag = [f"SEN_{idx}"] # Etiqueta única
28         all_docs.append(LabelDoc(words, tag))
```

1.1.3. Entrenamiento del modelo Doc2Vec:

- I. Se configura el modelo Doc2Vec con 50 dimensiones de vector, un mínimo de 2 ocurrencias de palabras, 10 épocas de entrenamiento y una tasa de aprendizaje que disminuye gradualmente.
- II. Se construye el vocabulario del modelo a partir de los documentos procesados.
- III. El modelo se entrena durante 10 épocas, ajustando los pesos con cada pasada y reduciendo la tasa de aprendizaje gradualmente.
- IV. El modelo entrenado se guarda en un archivo para su uso futuro.

```
39 # Entrenar el modelo Doc2Vec
40 model = Doc2Vec(vector_size=50, min_count=2, epochs=10, alpha=0.025, min_alpha=0.025)
41 model.build_vocab(all_docs)
42
43 print("Entrenando el modelo...")
44 for epoch in range(10):
45     print(f"Época {epoch+1}")
46     model.train(all_docs, total_examples=model.corpus_count, epochs=model.epochs)
47     model.alpha -= 0.002 # Reducir la tasa de aprendizaje
48     model.min_alpha = model.alpha
49
50 # Guardar el modelo
51 model.save('wikipedia_sample.doc2vec')
```

1.1.4. Evaluación del modelo:

- I. Se selecciona un documento aleatorio de la lista `all_docs` y se imprime su contenido.
- II. Se calcula la similitud entre el documento seleccionado y los otros documentos del modelo usando **`most_similar`**. Los 10 documentos más similares se recuperan y se imprimen junto con su nivel de similitud.

```

44 # Análisis de similitud
45 doc_id = np.random.randint(len(all_docs)) # Seleccionar un documento aleatorio
46 print(f"Documento objetivo (ID={doc_id}): {all_docs[doc_id].words}")
47
48 sims = model.dv.most_similar(f"SEN_{doc_id}", topn=10) # Documentos más similares
49 print("\nDocumentos similares:")
50 for sim_id, similarity in sims:
51     similar_idx = int(sim_id.replace("SEN_", ""))
52     print(f"r\"{sim_id} (similaridad: {similarity:.4f}): {all_docs[similar_idx].words}\"")

```

1.2 Resultado

- I. El código seleccionará una muestra aleatoria de hasta 5000 líneas del archivo wikipedia.txt. Si el archivo tiene menos de 5000 líneas, seleccionará todas las disponibles.

```
PS D:\Unidad7\2024-1\InteligenciaArtificial\Final\DataBoletines> python wikipedia_gensim_doc2vec.py
[total de documentos procesados: 4956]

Primeros 3 documentos: [LabelDoc(words=['leunos,', 'los,', 'leunos', 'en', 'tribin', 'leuni', 'eran', 'una', 'de', 'la', 'provincia', 'reana', 'de', 'vallencia', 'es', 'una', 'atlat', 'galicia', 'tendiendo', 'una', 'palmento', 'es', 'el', 'norte', 'de', 'portugal', 'entre', 'los', 'rios', 'lima', 'y', 'yafo', 'su', 'no', 'de', 'europoeo', 'leuks', 'no', 'que', 'significa', 'literalmente', 'brillante', 'resplandeciente', 'cfr', 'abelDoc(words=[club', 'basquet', 'lliria', 'el', 'club', 'basquet', 'lliria', 'es', 'un', 'equipo', 'da', 'la', 'ciudad', 'lliria', 'valencia', 'espaa', 'que', 'actualmente', 'juega', 'en', 'la', 'liga', 'eba', 'esta', 'est', 'espaol', 'el', 'baloncesto', 'nace', 'en', 'lliria', 'en', '1945', 'importado', 'desde', 'madrid', 'elizaiba', 'alli', 'el', 'servicio', 'militar', 'desnde', 'el', 'ascenso', 'en', 'la', 'temporada', '67', 'fueron', 'vencidos', 'por', 'los', 'viros', 'propos', 'gracias', 'testos', 'adquisia', 'temas', 'entre', 'la', '2', 'y', 'la', '1', 'division', 'nacional', 'hasta', 'que', 'finalmente', 'en', '1982', 'se', 'harise', 'subcampeon', 'de', 'la', '1', 'primera', 'division', 'y', 'consegue', 'el', 'ascenso', 'a', 'la', '1', 'division', 'tecnico', 'andreu', 'casadevall', 'competicon', 'en', 'la', 'que', 'permanecerla', 'durante', 'las', 'ras', 'jugar', 'las', 'ultimas', 'temporadas', 'en', 'la', 'liga', 'eba', 'el', 'actual', 'club', 'ha', 'rada', '2005', '2006', 'en', 'la', 'leb', '2', 'en', 'la', 'que', 'consigulo', 'la', 'permanencia', 'en', 'la', '1', 'division', 'nacional', 'como', 'personas', 'importantes', 'que', 'han', 'pasado', 'en', 'este', 'a', 'ser', 'de', 'los', 'fundadores', 'del', 'club', 'de', 'baloncesto', 'recuerdo', 'sea', 'de', 'los', 'seales', 'de', 'los', '80', 'formada', 'por', 'dan', 'palombizio', 'palomo', 'y', 'vernos', 'smith', 'por', 'el', '2001', 'se', 'fundo', 'el', 'club', 'esportiu', 'lliria', 'para', 'continuar', 'la', 'tradicion', 'escudo', 'del', 'lliria', 'representa', 'un', 'balon', 'de', 'baloncesto', 'amarillo', 'con', 'unos', 'colores', 'del', 'club', 'con', 'los', 'bordes', 'amarillos', 'y', 'negros', 'introduciendo', 'en', 'u', 'vez', 'la', 'palabra', 'lliria', 'en', 'minusculas', 'el', 'uniforme', 'es', 'predominantemente
```

- II. El modelo Doc2Vec se entrenará en 10 épocas con las oraciones preprocesadas. En cada época, el modelo ajustará sus parámetros para mejorar la representación de los documentos.

- III. Durante el entrenamiento, la tasa de aprendizaje (α) disminuirá gradualmente en 0.002 después de cada época.

```
la', 'especie', 'mas', 'caracteristica', 'del', 'valle', 'c  
n', 'paisaje', 'de', 'singular', 'belleza'], tags=['SEN_2']  
Entrenando época 1  
Entrenando época 2  
Entrenando época 3  
Entrenando época 4  
Entrenando época 5  
Entrenando época 6  
Entrenando época 7  
Entrenando época 8  
Entrenando época 9  
Entrenando época 10
```

- IV. Se muestra el contenido del documento seleccionado aleatoriamente como punto de partida del análisis.

```
Documento objetivo (ID=2911): ['comendador', 'municipio', 'comendador', 'es', 'el', 'mu',  
'piña', 'su', 'nombre', 'honra', 'a', 'nicolas', 'de', 'ovando', 'comendador', 'de',  
'vivienda', 'de', '2002', 'el', 'municipio', 'tiene', 'una', 'poblacion', 'total', 'de',  
'mujeres', 'y', '12901', 'mujeres', 'la', 'poblacion', 'urbana', 'del', 'municipio', 'era',  
'nacionales', 'incluyen', 'la', 'poblacion', 'de', 'los', 'distritos', 'municipales', 'gu',  
'embre', 'de', '1930', 'se', 'sustituyo', 'el', 'nombre', 'de', 'comendador', 'para', 'e',  
'piña', 'la', 'ley', 'no', '342', 'del', '29', 'de', 'mayo', 'de', '1972', 'restituyo',  
, 'fue', 'fundada', 'comendador', 'la', 'principal', 'actividad', 'economica', 'del', 'com',  
'servicios', 'terciarios', 'han', 'ido', 'incrementandose', 'sobre', 'todo', 'el', 'comer
```

- V. Se calcularán los 10 documentos más similares según el modelo, como el que se presenta a continuación:

```
SEN_4770 (similaridad: 0.5726): ['yohualtecutli', 'yohualtecutli', 'en', 'nahuatl', 'se',  
cano', 'de', 'la', 'noche', 'y', 'protegia', 'el', 'sueño', 'de', 'los', 'niños', 'junto',  
n', 'descender', 'de', 'ancestros', 'comunes', 'del', 'panton', 'mesoamericano', 'los',  
li', 'como', 'un', 'señor', 'de', 'las', 'horas', 'nocturnas', 'y', 'yacahuitzli', 'era',  
'la', 'puesta', 'del', 'sol', 'las', 'fuentes', 'clasicas', 'lo', 'asocian', 'por', 'e',  
la', 'incineracion', 'las', 'conclusiones', 'del', 'periodo', 'y', 'la', 'colocacion', 'd',  
rono', 'se', 'considero', 'que', 'seria', 'una', 'de', 'las', 'tres', 'piedras', 'del', 'd',  
'centro', 'de', 'las', 'casas', 'mayas', 'que', 'era', 'una', 'imagen', 'a', 'escala', 'd',  
'considera', 'que', 'las', 'estrellas', 'localizadas', 'bajo', 'el', 'cinturon', 'de', 'o',  
y', 'yacahuitzli', 'se', 'identifican', 'en', 'esta', 'region', 'con', 'orion', 'ubicado',  
ferencia', 'en', 'la', 'universidad', 'de', 'leiden']
```

2. Entrenando al modelo gensim_doc2vec con conjunto de datos denominado boletines (paso realizados, entrenamiento con la data boletines y printscreen de corrida...)

2.1. Pasos realizados

2.1.1. Carga de datos:

- I. Se utilizó el archivo boletines.csv, compuesto por los campos Idtexto, Texto y Supervision, separados por punto y coma (;).
- II. Se concatenó el contenido del campo Texto en un solo bloque para su procesamiento posterior.

```
8 # Leer el archivo CSV  
9 with open("boletines.csv", newline='', encoding='UTF-8', errors='ignore') as csvfile:  
10     reader = csv.reader(csvfile)  
11     data = ''  
12     count = 0  
13     for row in reader:  
14         count += 1  
15         if count > 55:  
16             break  
17         # Verificar que la fila tenga al menos dos columnas  
18         if len(row) > 1:  
19             data += row[1] # Concatenar solo la columna de texto  
20         else:  
21             print(f"fila {count} no tiene suficientes columnas: {row}")
```

2.1.2. Preprocesamiento de datos:

- I. Se dividieron los textos en oraciones utilizando una expresión regular.
- II. Se eliminaron caracteres de puntuación para limpiar los datos.

```
23 # Configurar una expresión regular para dividir el párrafo en oraciones.  
24 sentenceEnders = re.compile('[.?!']*)  
25 data_list = sentenceEnders.split(data)  
26  
27 # Crear un namedtuple con words y tags  
28 LabelDoc = namedtuple('LabelDoc', 'words tags')  
29 exclude = set(string.punctuation)  
30 all_docs = []  
31 count = 0  
32 for sen in data_list:  
33     word_list = sen.split()  
34     if len(word_list) < 3:  
35         continue  
36     tag = ['SEN_' + str(count)]  
37     count += 1  
38     sen = ''.join(ch for ch in sen if ch not in exclude)  
39     all_docs.append(LabelDoc(sen.split(), tag))  
40  
41 print(all_docs[0:10])
```

- III. Se descartaron oraciones con menos de tres palabras, ya que aportan poca información semántica.
- IV. Cada oración se etiquetó de manera única (SENT_0, SENT_1, entre otros).

2.1.3. Entrenamiento del modelo Doc2Vec:

- I. Se inicializa un modelo Doc2Vec con parámetros específicos (tasa de aprendizaje fija).
- II. Se construye un vocabulario a partir de los documentos.
- III. El modelo se entrena en 10 iteraciones, ajustando los pesos para optimizar los vectores de documentos.
- IV. Finalmente, se guarda el modelo entrenado para usarlo posteriormente en tareas como inferencia o similitud de documentos.

```

43 # Entrenar el modelo Doc2Vec
44 model = Doc2Vec(alpha=0.025, min_alpha=0.025) # usar tasa de aprendizaje fija
45 model.build_vocab(all_docs)
46 for epoch in range(10):
47     model.train(all_docs, total_examples=model.corpus_count, epochs=model.epochs)
48     model.alpha -= 0.002 # disminuir la tasa de aprendizaje
49     model.min_alpha = model.alpha # fijar la tasa de aprendizaje, sin decaimiento
50
51 model.save('boletinesModel.doc2vec')

```

2.1.4. Evaluación del modelo:

- I. Se selecciona un documento objetivo aleatorio y se imprime
- II. Se utiliza el modelo Doc2Vec para encontrar los documentos más cercanos en el espacio vectorial.
- III. Se imprimen los resultados:
Mostrando el contenido del documento objetivo y hasta 8 documentos similares, junto con sus identificadores.

```

54 # Obtener un documento aleatorio
55 doc_id = np.random.randint(model.dv.vectors.shape[0])
56 print(doc_id)
57
58 sims = model.dv.most_similar(doc_id, topn=model.dv.vectors.shape[0])
59 print('TARGET', all_docs[doc_id].words)
60
61 count = 0
62 for i in sims:
63     if count > 8:
64         break
65     pid = int(i[0].replace("SEN_", ""))
66     print(i[0], ":", all_docs[pid].words)
67     count += 1

```

2.2. Resultados:

- I. **LabelDoc(words)**: Representa el contenido del texto del documento después de separar palabras y eliminar signos de puntuación.

```
PS D:\Universidad\2024-1\InteligenciaArtificial\Final\DataBoletines> python boletines_gensim_doc2vec.py
[LabelDoc(words=[Textopagboja, boletín, oficial, de, la, junta, de, andalucía, número, 1, de, diciembre, de, 2016, página, 134, 5], tags=[SEN_0]), LabelDoc(words=[otros, ciales, consejería, de, salud, anuncio, de, 25, de, noviembre, de, 2016, de, la, territorial, de, igualdad, salud, y, políticas, sociales, en, almería, para, la, e, edicto, del, citado, acto], tags=[SEN_1]), LabelDoc(words=[de, conformidad, con, l, en, los, artículos, 59], tags=[SEN_2]), LabelDoc(words=[5, y, 61, de, la, ley, 36, de, noviembre, de, régimen, jurídico, de, las, administraciones, públicas, y, del, administrativo, común, al, no, haberle, podido, ser, practicada, notificación, a, jesús, rojas, sáez, y, doña, maría, ángeles, canón, frías, se, pública, extracto dictado], tags=[SEN_3]), LabelDoc(words=[para, su, conocimiento, integro, podrán, compare servicio, de, protección, de, menores, de, la, delegación, territorial, de, igual, políticas, sociales, de, almería, sito, en, carretera, de, ronda, edif], tags=[SEN_4]), LabelDoc(words=[SEN_5]), LabelDoc(words=[SEN_6]), LabelDoc(words=[SEN_7]), LabelDoc(words=[de, conformidad, con, lo, dispuesto, en, el, art], tags=[SEN_7]), LabelDoc(words=[boletín, oficial, de, la, junta, de, andalucía, número, 233, lunes, 5, de, 2016, página, 1], tags=[SEN_8]), LabelDoc(words=[disposiciones, generales, consejería, de, vivienda, resolución, de, 28, de, noviembre, de, 2016, de, la, dirección, general, por, la, que, se, convocan, pruebas, para, la, obtención, del, certificado, profesional, acreditativo, de, la, cualificación, inicial, de, los, conductores, de, vehículos, destinados, al, transporte, por, carretera, y, se, detalla, la, composición, tribunales, calificadoros, así, como, las, fechas, horarios, y, lugares, de, cal, las, pruebas, en, el, año, 2017], tags=[SEN_9])]
```

- II. Muestra el contenido del documento seleccionado aleatoriamente como punto de partida del análisis.

```
TARGET ['5, y, 61, de, la, ley, 301992, de, 26, de, noviembre, de, régimen, jurídico, de, l, as, administraciones, públicas, y, del, procedimiento, administrativo, común, modificada, por, la, ley, 41999, de, 13, de, enero, se, la, anuncia, que, por, el, director, general, de, infraestructuras, de, la, consejería, de, fomento, y, vivienda, se, ha, dictado, resolución, sancionada es, en, al, procedimiento, administrativo, de, expedite, sancionador, en, materia, de, carreteras, con, referencia, 200716, por, el, incumplimiento, depagboja, boletín, oficial, de, la, junta, de, andalucía, número, 30, martes, 14, de, febrero, de, 2017, página, 140, 5], SEN_54 : ['5, de, la, ley, 3092, de, 26, de, noviembre, de, régimen, jurídico, de, las, administraciones, públicas, y, del, procedimiento, administrativo, común, se, notifica, la, ipagboja, boletín, oficial, de, la, junta, de, andalucía, número, 10, martes, 17, de, enero, de, 2017, página, 134, 5]
```

- III. Se muestra el contenido similar: Representado como una lista de palabras preprocesadas.

```
SEN_147 (similitud: 0.9792) : ['b, dependencia, que, tramita, el, expediente, dirección, general, de, infraestructuras']
SEN_38 (similitud: 0.9719) : ['b, dependencia, que, tramita, el, expediente, dirección, general, de, infraestructuras']
SEN_122 (similitud: 0.9658) : ['b, dependencia, que, tramita, el, expediente, dirección, general, de, infraestructuras']
SEN_320 (similitud: 0.9388) : ['b, dependencia, que, tramita, el, expediente, secretaría, general, técnica']
SEN_237 (similitud: 0.9319) : ['b, dependencia, que, tramita, el, expediente, secretaría, general, técnica']
SEN_222 (similitud: 0.9237) : ['b, dependencia, que, tramita, el, expediente, secretaría, general, técnica']
SEN_327 (similitud: 0.9119) : ['b, dependencia, que, tramita, el, expediente, secretaría, general, técnica']
SEN_67 (similitud: 0.9084) : ['b, dependencia, que, tramita, el, expediente, dirección, general, de, movilidad']
SEN_380 (similitud: 0.9075) : ['expediente, inaga, 50020101201610185']
```

SEN: Es el identificador (tag) del documento similar. Los identificadores se asignan al preprocesar el corpus (por ejemplo, SEN_147, SEN_38, SEN_122, entre otros).

Similitud: Representa el porcentaje de similitud que tiene con el objetivo, el SEN 147 posee una similitud de 0.9792, lo cual es cercano y superior al siguiente SEN_38, que posee una similitud de 0.9719.

Relaciones semánticas captadas: El modelo ha detectado relaciones significativas entre conceptos de física relacionados con la gravedad.

Jerarquía de similitud: Los documentos con mayor contenido compartido (términos y contexto) reciben puntajes de similitud más altos.

3. Entrenando al modelo gensim_doc2vec con conjunto de datos denominado boletines de prueba (paso realizados, entrenamiento con boletines de prueba y printscreen de corrida...)

3.1. Pasos realizados

3.1.1. Carga de datos:

Se utilizó el archivo **boletinesprueba.csv**, que posee los campos Idtexto, Texto y Supervisión, pero a diferencia de **boletines.csv** posee los valores separados por coma (.). Por lo que cambiamos el delimitador.

Se repiten los mismos pasos que con boletines.csv: **2.1.2. Preprocesamiento de datos,**

2.1.3. Entrenamiento del modelo Doc2Vec, 2.1.4. Evaluación del modelo.

3.2. Resultados

- Se selecciona un documento aleatorio del corpus (TARGET) para compararlo con el resto. (Este documento está representado como una lista de palabras procesadas (sin puntuación))
- El modelo Doc2Vec busca los documentos más similares al objetivo en el espacio semántico. La similitud se mide como un valor entre -1 (completamente opuesto) y 1 (idéntico).
- Los resultados incluyen el identificador del documento (SEN_X), su contenido preprocesado y el puntaje de similitud.

Interpretación de los Resultados:

- Los documentos con puntajes de similitud más altos comparten términos clave o temas relacionados con el documento objetivo.

- Los documentos menos similares también pueden estar relacionados de manera más vaga, pero tienen menos coincidencias directas en palabras o contexto.

```

Windows PowerShell
PS D:\Universidad\2024-1\InteligenciaArtificial\Final\DataBoletines> python boletinesprueba_gensim_doc2vec.py

[LabelDoc(words=['contratación', 'del', 'sector', 'público', 'ministerio', 'de', 'fomento', '66033', 'anuncio', 'de', 'licitación', 'de', 'dirección', 'general', 'de', 'carreteras'], tags=['SEN_0']), LabelDoc(words=['objeto', 'construcción', 'y', 'explotación', 'del', 'área', 'de', 'servicio', 'de', 'cieza', 'en', 'la', 'auto', 'vía', 'a67', 'cantabriamesetaadministración', 'de', 'justicia', '52917', 'juzgado', 'de', 'lo', 'social', 'número', '4', 'de', 'almería', 'procedimiento', 'procedimiento', 'ordinario', '12222014', 'negociado', 'cn', 'n'], tags=['SEN_1']), LabelDoc(words=['0401344s20140004969', 'de', 'dd'], tags=['SEN_2']), LabelDoc(words=['david', 'lopez', 'asensio', 'abogado', 'contra', 'dd'], tags=['SEN_3']), LabelDoc(words=['participaciones', 'societarias', 'tecnológicas', 'y', 'empeariales', 'sl'], tags=['SEN_4']), LabelDoc(words=['madar', 'desarroll', 'os', 'sl'], tags=['SEN_5']), LabelDoc(words=['exclusive', 'marketing', 'hotels', 'tourism', 'leisure', 's'], tags=['SEN_6']), LabelDoc(words=['y', 'cymar', 'gestion', 'hotelera', 's'], tags=['SEN_7']), LabelDoc(words=['abogado', 'e', 'd', 'i', 'c', 't', 'o', 'cédula', 'de', 'citación', 'en', 'resolución', 'del', 'día', 'de', 'la', 'fecha', 'dictada', 'en', 'la', 'ejecución', 'número'], tags=['SEN_8']), LabelDoc(words=['seguidos', 'en', 'este', 'juzgado', 'de', 'lo', 'social', 'número', '4', 'de', 'almería', 'y', 'su', 'provincia', 'en', 'mat', 'eria', 'de', 'procedimiento', 'ordinarioadministració', 'local', 'diputació', 'de', 'barcelona'], tags=['SEN_9'])]
13
TARGET ['presidenta', 'de', 'la', 'diputació', 'de', 'barcelona', 'resta', 'aprovat', 'definitivament', 'el', 'projecte', 'constructiu', 'ampliació', 'del', 'pont', 'de', 'can', 'molas', 'i', 'millora', 'de', 'laccés', 'del', 'veinat', 'de', 'vista', 'alegre', 'a', 'la', 'carretera', 'b522']
SEN_15 (similaridad: 0.9966) : ['per', 'impugnar', 'aquesta', 'resolució', 'que', 'posa', 'fi', 'a', 'la', 'via', 'administrativa', 'cal', 'interposar', 'recurs', 'contenciós', 'a1', 'diciembre', '2016', 'boletín', 'oficial', 'de', 'la', 'provincia', 'de', 'huesca', 'nº', '230', 'administración', 'del', 'estado', 'ministerio', 's', 'ministerio', 'de', 'fomento', 'dirección', 'general', 'de', 'carreteras', 'demarcación', 'de', 'carreteras', 'del', 'estado', 'en', 'aragón', 'anuncio', '5143', 'unidad', 'de', 'carreteras', 'huesca', 'por', 'el', 'servicio', 'de', 'obras', 'públicas', 'y', 'patrimonio', 'de', 'la', 'diputación', 'provincial', 'de', 'huesca']
SEN_21 (similaridad: 0.9964) : ['m1', 'en', 'cumplimiento', 'de', 'lo', 'preceptuado', 'en', 'el', 'artículo', '83', 'de', 'la', 'ley', '392015', 'de', '1', 'de', 'octubre', 'del', 'procedimiento', 'administrativo', 'común', 'de', 'las', 'administraciones', 'públicas', 'y', 'de', 'conformidad', 'con', 'lo', 'dispuesto', 'en', 'los', 'artículos', '18', 'y', '19']
SEN_29 (similaridad: 0.9956) : ['83170', 'euros', 'con', 'las', 'siguientes', 'prescripciones', 'a', 'cumplimentar', 'en', 'el', 'proyecto', 'de', 'construcción', '1']
SEN_26 (similaridad: 0.9951) : ['la', 'dirección', 'general', 'de', 'carreteras', 'con', 'fecha', '30', 'de', 'diciembre', '2016', 'en', 'uso', 'de', 'las', 'competencias', 'establecidas', 'en', 'el', 'artículo', '32', 'del', 'vigente', 'reglamento', 'general', 'de', 'carreteras', 'de', '2', 'de', 'septiembre', 'de', '1994', 'ha', 'resuelto', '1']
SEN_12 (similaridad: 0.9950) : ['2', 'del', 'reglament', 'dobres', 'activitats', 'i', 'serveis', 'dels', 'en', 's', 'locals', 'aprovat', 'pel', 'decret', '1791995', 'de', '13', 'de', 'juny', 'es', 'fa', 'saben', 'que', 'en', 'virtut', 'de', 'resolució', 'dictada', 'en', 'data', '14', 'doctubre', 'de', '2016', 'per', 'lexcma']
SEN_19 (similaridad: 0.9950) : ['aprobado', 'por', 'la', 'dirección', 'general', 'de', 'carreteras', 'por', 'resolución', 'de', '14', 'de', 'abril', 'de', '2016', 'el', 'proyecto', 'de', 'construcción', 'de', 'autovía', 'gr43', 'tramo', 'pinos', 'puente', '-', 'atarfe', 'enlace', 'con', 'la', 'futura', 'a44']
SEN_23 (similaridad: 0.9949) : ['otros', 'anuncios', 'oficiales', 'ministerio', 'de', 'fomento', '2044', 'anuncio', 'de', 'la', 'demarcación', 'de', 'carreteras', 'del', 'estado', 'en', 'la', 'rioja', 'por', 'el', 'que', 'se', 'publica', 'la', 'aprobación', 'provisional', 'y', 'se', 'ordena', 'la', 'incoación', 'del', 'expediente', 'de', 'información', 'pública', 'exclusivamente', 'a', 'los', 'efectos', 'de', 'la', 'ley', 'de', '16', 'de', 'diciembre', 'de', '1954', 'de', 'expropiación', 'forzosa', 'del', 'proyecto', 'de', 'trazado', 't41', 'o5540', 'autovía', 'a68']
SEN_10 (similaridad: 0.9948) : ['gerència', 'de', 'serveis', 'dinfraestructures', 'viàries', 'i', 'mobilitat', 'anunci', 'sobre', 'aprovació', 'definitiva', 'del', 'projecte', 'constructiu', 'ampliació', 'del', 'pont', 'de', 'can', 'molas', 'i', 'millora', 'de', 'laccés', 'del', 'veinat', 'de', 'vista', 'alegre', 'a', 'la', 'carretera', 'b522']
SEN_14 (similaridad: 0.9946) : ['manlleu', 'en', 'no', 'haber', 'estat', 'formulades', 'reclamacions', 'o',

```


Conclusiones

Aplicabilidad de Doc2Vec: El modelo Doc2Vec, basado en Gensim, demuestra ser una herramienta efectiva para representar documentos como vectores densos, lo que facilita tareas como la clasificación, búsqueda y agrupamiento de textos. Su capacidad para capturar el contexto semántico global de los textos lo hace especialmente útil para conjuntos de datos especializados, como los boletines analizados.

Versatilidad del entrenamiento: El uso de datasets variados, como Wikipedia.txt y los archivos de boletines, evidencia la flexibilidad de Doc2Vec para adaptarse a diferentes dominios y contextos de trabajo. Entrenar el modelo con datos específicos mejora significativamente su capacidad para interpretar y analizar textos similares.

Importancia de la preprocesamiento: La calidad de los resultados obtenidos depende en gran medida del preprocesamiento de los datos. Pasos como la limpieza de texto, eliminación de caracteres irrelevantes y segmentación adecuada de frases son críticos para garantizar un entrenamiento efectivo del modelo.

Gensim como herramienta clave: Gensim se posiciona como una biblioteca robusta y accesible para el procesamiento de lenguaje natural, ofreciendo una implementación eficiente de Doc2Vec que permite a los investigadores entrenar modelos personalizados con relativa facilidad.

Perspectivas futuras: El enfoque empleado puede extenderse a otros conjuntos de datos y adaptarse para mejorar la comprensión semántica en diversas áreas, como análisis de sentimientos, recuperación de información y generación de resúmenes automáticos.

Referencias Bibliográficas

Le, Q., & Mikolov, T. (2014). Representaciones Distribuidas de Frases y Documentos.
Goldberg, Y., & Levy, O. (2014). Word2Vec Explicado: Derivando el Método de Embedding de Palabras con Muestreo Negativo de Mikolov.

Referencias Web

<https://radimrehurek.com/gensim/>: Radim Řehůřek. (2018). Gensim: Modelado de Temas para Humanos.

<https://radimrehurek.com/gensim/models/doc2vec.html>: Documentación de Doc2Vec en Gensim.

<https://machinelearningmastery.com>: Mastery en Aprendizaje Automático. (2021). Cómo Desarrollar un Modelo de Vector de Párrafo.

https://es.wikipedia.org/wiki/Aprendizaje_profundo: Colaboradores de Wikipedia. (2023). Aprendizaje Profundo.

<https://towardsdatascience.com>: Hacia la Ciencia de Datos. (2022). Doc2Vec y Word2Vec: Diferencias y Aplicaciones.