

Anexo 4

Proyecto 4: Creación de Words Embedding Nivel 1

Mg. Luis Felipe Bustamante Narváez

Word2Vec

Es un modelo que se utiliza para aprender representaciones vectoriales de palabras. Estas representaciones pueden capturar muchas propiedades lingüísticas de las palabras, como su significado semántico, gramatical y hasta contextual.

Para este primer ejemplo, usaremos un texto corto llamado mundiales que tiene alrededor de 50.000 caracteres.

Librerías

```
In [...] !pip install pypdf2
```

```
In [...] import string
from gensim.models import Word2Vec
import PyPDF2
```

Cargamos el documento

```
In [...] with open('Entrenamiento_Word2Vec/mundiales.txt',
                'r', encoding='utf-8') as file:
    documento = file.read()
```

```
In [...] documento
```

```
In [...] len(documento)
```

```
Out[... 48155
```

Procesamiento de datos

El objetivo del procesamiento es convertir el documento en una lista de frases, y cada frase en una lista de palabras, eliminando signos de puntuación y convirtiendo todo a minúsculas.

```
In [...] # Dividimos el documento en frases usando la coma como separador
frases = documento.split(',')
len(frases)
```

```
Out[... 533
```

```
In [...] # Mostramos un ejemplo
frases[0]
```

```
Out[... 'Capítulo 1: Historia del Mundial de Fútbol\nLos primeros pasos del Mundial\nLa historia d
el Mundial de Fútbol se remonta a principios del siglo XX'
```

```
In [... # Mostramos un ejemplo
frases[500]
```

```
Out[... ' **Adidas**'
```

```
In [... # Limpiamos las frases
frases_limpias = []
for frase in frases:
    #Eliminamos la puntuación y dividimos por espacios
    tokens = frase.translate(str.maketrans('', '', string.punctuation)).split()
    #print(tokens) #para mostrar qué ha hecho hasta aquí
    #Convertimos a minúsculas
    tokens = [word.lower() for word in tokens]
    #print(tokens) #para mostrar qué ha hecho hasta aquí
    if tokens:
        frases_limpias.append(tokens)
```

```
In [... # Mostramos los resultados
frases_limpias[500]
```

```
Out[... ['adidas']
```

Entrenamiento del modelo Word2Vec

```
In [... model = Word2Vec(sentences=frases_limpias,
                        vector_size=500,
                        window=5,
                        min_count=1,
                        workers=4)
```

Explicación:

- sentences: Es la lista de palabras que vamos a vectorizar
- vector_size: Es el tamaño de dimensiones que le daremos al vector
- window: Son la cantidad de palabras por encima y por debajo que le darán contexto
- min_count: La aparición mínima de una palabra para tenerla en cuenta en el entrenamiento
- workers: Cantidad de núcleo de procesador que vamos a invertir en el entrenamiento

Pruebas

```
In [... # Verificamos el vector para alguna palabra
vector = model.wv['mundial']
vector
```

```
In [... # Mostramos las palabras cercanas
palabras_cercanas = model.wv.most_similar('jugador', topn=10)
palabras_cercanas
# Es probable que la similitud falle por tener tan pocas palabras en el texto
```

```
Out[... [('historia', 0.9089663624763489),
        ('su', 0.9083009362220764),
        ('con', 0.9079578518867493),
        ('y', 0.9078695774078369),
        ('brasil', 0.9078560471534729),
        ('en', 0.9078254699707031),
        ('más', 0.9077353477478027),
        ('del', 0.9077184200286865),
        ('francia', 0.9076839685440063),
        ('juego', 0.9075537919998169)]
```

Guardar modelo

```
In [... model.save('Entrenamiento_Word2Vec/mundiales.model')
```

Cargar el modelo

```
In [... modelo_cargado = Word2Vec.load('Entrenamiento_Word2Vec/mundiales.model')
```

```
In [... # Probamos con el modelo caragado
palabras_cercanas2 = modelo_cargado.wv.most_similar('jugador', topn=10)
palabras_cercanas2
```

```
Out[... [('historia', 0.9089663624763489),
        ('su', 0.9083009362220764),
        ('con', 0.9079578518867493),
        ('y', 0.9078695774078369),
        ('brasil', 0.9078560471534729),
        ('en', 0.9078254699707031),
        ('más', 0.9077353477478027),
        ('del', 0.9077184200286865),
        ('francia', 0.9076839685440063),
        ('juego', 0.9075537919998169)]
```

Guardar Embedido

Existen dos maneras, usando .txt sin binarios, y usando .bin con binarios.

```
In [... model.wv.save_word2vec_format('Entrenamiento_Word2Vec/mundiales_emb.txt', binary=False)
model.wv.save_word2vec_format('Entrenamiento_Word2Vec/mundiales_emb.bin', binary=True)
```

Cargar Embedidos

Si se carga el .txt, se usa sin binarios; si se carga el .bin, se usa con binarios

```
In [... from gensim.models import KeyedVectors
embedding_cargado_txt = KeyedVectors.load_word2vec_format(
    'Entrenamiento_Word2Vec/mundiales_emb.txt', binary=False)
```

```
In [... embedding_cargado_bin = KeyedVectors.load_word2vec_format(
    'Entrenamiento_Word2Vec/mundiales_emb.bin', binary=True)
```

```
In [... # Probamos
```

```
embedding_cargado_txt
```

```
Out[... <gensim.models.keyedvectors.KeyedVectors at 0x21302849a30>
```

```
In [... # Probamos  
embedding_cargado_bin
```

```
Out[... <gensim.models.keyedvectors.KeyedVectors at 0x213062ccaa0>
```

Analogías

```
In [... def analogics(v1, v2, v3):  
        simil = embedding_cargado_bin.most_similar(positive=[v1,v3],  
                                                    negative=[v2]  
                                                    )  
        print(f'{v1} es a {v2}, como {simil[0][0]} es a {v3}')
```

```
In [... analogics('jugador', 'fútbol', 'historia')  
jugador es a fútbol, como brasil es a historia
```

Conclusiones

El texto mundiales tiene cerca de 50.000 caracteres, lo que implica una base de datos pequeña para entrenar un modelo. De cierta forma, el modelo se ajusta en algunos casos puntuales, pero suele mostrar demasiadas stopwords, que tendríamos que manipular para mejorar la predicción de analogías. Veremos con un texto más grande, como se generaría la predicción, por ejemplo el libro "Cien años de soledad" de Gabriel García Márquez.

Mg. Luis Felipe Bustamante Narváez