

Tarea 10

NLP

Descripción: Dados 3 textos construir la matriz término-documento a partir de bigramas de palabras y trigramas de caracteres y determinar la similitud coseno.

Alumno: Miguel Angel Soto Hernandez

```
1 import nltk
2 import urllib.request as urllib2
3 from bs4 import BeautifulSoup
4 nltk.download('book')
5 from nltk.book import *
6 import matplotlib.pyplot as plt
7 import numpy as np
8 import pandas as pd
9 import plotly.express as px
```

```
1 # Moby Dick
2 texto_1 = text1
3
4 # Sense and Sensibility
5 texto_2 = text2
6
7 # The book of genesis
8 texto_3 = text3
```

```
1 import re
2 def normalizar_texto(texto, quitar_numeros, incluir_puntuacion, convertir_a_minusculas, incluir_stopwords):
3     if quitar_numeros == True:
4         texto = [re.sub('\d', ' ', palabra) for palabra in texto]
5
6     if incluir_puntuacion == True:
7         texto = [palabra for palabra in texto if re.search('\w+', palabra)]
8
9     if convertir_a_minusculas == True:
10        texto = [palabra.lower() for palabra in texto]
11
12    if incluir_stopwords == True:
13        stopwords = sorted(nltk.corpus.stopwords.words('english'))
14        texto = [palabra for palabra in texto if palabra not in stopwords]
15
16    return texto
```

```
1 texto_1 = normalizar_texto(texto_1, quitar_numeros=True,
2                             incluir_puntuacion=True,
3                             convertir_a_minusculas=True,
4                             incluir_stopwords=True)
5 print(texto_1[:10])
```

```
['moby', 'dick', 'herman', 'melville', 'etymology', 'supplied', 'late', 'consumptive', 'usher']
```

```

1 texto_2 = normalizar_texto(texto_2, quitar_numeros=True,
2                             incluir_puntuacion=True,
3                             convertir_a_minusculas= True,
4                             incluir_stopwords=True)
5 print(texto_2[:10])

```

```
['sense', 'sensibility', 'jane', 'austen', 'chapter', 'family', 'dashwood', 'long', 'settled',
```

```

1 texto_3 = normalizar_texto(texto_3, quitar_numeros=True,
2                             incluir_puntuacion=True,
3                             convertir_a_minusculas= True,
4                             incluir_stopwords=True)
5 print(texto_3[:10])

```

```
['beginning', 'god', 'created', 'heaven', 'earth', 'earth', 'without', 'form', 'void', 'darknes
```

```

1 def obtener_bigramas(texto):
2     bigramas = list(nltk.bigrams(texto))
3     return bigramas
4
5
6 def contar_bigramas(bigramas):
7     frecuencia_bigramas = FreqDist(bigramas)
8     return frecuencia_bigramas

```

```

1 bigramas_texto_1 = obtener_bigramas(texto_1)
2 print(bigramas_texto_1[:10])
3
4 frecuencia_bigramas_texto_1 = contar_bigramas(bigramas_texto_1)
5 print(frecuencia_bigramas_texto_1[:10])

```

```
[(('moby', 'dick'), ('dick', 'herman'), ('herman', 'melville'), ('melville', 'etymology'), ('etymology', 'man')),
((('sperm', 'whale'), 182), (('white', 'whale'), 106), (('moby', 'dick'), 84), (('old', 'man'), 40))]
```

```

1 bigramas_texto_2 = obtener_bigramas(texto_2)
2 print(bigramas_texto_2[:10])
3
4 frecuencia_bigramas_texto_2 = contar_bigramas(bigramas_texto_2)
5 print(frecuencia_bigramas_texto_2[:10])

```

```
[(('sense', 'sensibility'), ('sensibility', 'jane'), ('jane', 'austen'), ('austen', 'chapter'), ('chapter', 'family'), ('family', 'dashwood'), ('dashwood', 'long'), ('long', 'settled'), ('settled', 'man')),
((('mrs', 'jennings'), 230), (('colonel', 'brandon'), 132), (('mrs', 'dashwood'), 121), (('sir', 'elton'), 40))]
```

```

1 bigramas_texto_3 = obtener_bigramas(texto_3)
2 print(bigramas_texto_3[:10])
3
4 frecuencia_bigramas_texto_3 = contar_bigramas(bigramas_texto_3)
5 print(frecuencia_bigramas_texto_3[:10])

```

```
[(('beginning', 'god'), ('god', 'created'), ('created', 'heaven'), ('heaven', 'earth'), ('earth', 'earth'), ('earth', 'without'), ('without', 'form'), ('form', 'void'), ('void', 'darkness'), ('darkness', 'man')),
((('said', 'unto'), 180), (('came', 'pass'), 63), (('thou', 'shalt'), 50), (('thou', 'hast'), 40))]
```

```

1 # Convertir palabras (tokens) en ngramas de caracteres
2 def ngramas_caracteres(texto, ngramas):
3     caracter_ngramas = [''.join(caracteres) for palabra in texto
4         for caracteres in zip(*[palabra[i:] for i in range(ngramas)])]
5     return caracter_ngramas
6
7
8 def frecuencia_ngramas_caracteres(ngramas_caracteres):
9     frecuencia_ngramas = FreqDist(ngramas_caracteres)
10    return frecuencia_ngramas

```

```

1 ngramas_caracteres_texto_1 = ngramas_caracteres(texto_1, 3)
2 print(ngramas_caracteres_texto_1[:10])
3
4 frecuencia_ngramas_caracteres_texto_1 = frecuencia_ngramas_caracteres(ngramas_caracteres_texto_1)
5 print(frecuencia_ngramas_caracteres_texto_1.most_common(10))

```

['mob', 'oby', 'dic', 'ick', 'her', 'erm', 'rma', 'man', 'mel', 'elv']
 [('ing', 7733), ('ter', 2292), ('ent', 2193), ('hal', 2135), ('wha', 1914), ('ale', 1903), ('or

```

1 ngramas_caracteres_texto_2 = ngramas_caracteres(texto_2, 3)
2 print(ngramas_caracteres_texto_2[:10])
3
4 frecuencia_ngramas_caracteres_texto_2 = frecuencia_ngramas_caracteres(ngramas_caracteres_texto_2)
5 print(frecuencia_ngramas_caracteres_texto_2.most_common(10))

```

['sen', 'ens', 'nse', 'sen', 'ens', 'nsi', 'sib', 'ibi', 'bil', 'ili']
 [('ing', 3355), ('ion', 1794), ('ent', 1726), ('eve', 1464), ('tio', 1418), ('ter', 1304), ('ve

```

1 ngramas_caracteres_texto_3 = ngramas_caracteres(texto_3, 3)
2 print(ngramas_caracteres_texto_3[:10])
3
4 frecuencia_ngramas_caracteres_texto_3 = frecuencia_ngramas_caracteres(ngramas_caracteres_texto_3)
5 print(frecuencia_ngramas_caracteres_texto_3.most_common(10))

```

['beg', 'egi', 'gin', 'inn', 'nni', 'nin', 'ing', 'god', 'cre', 'rea']
 [('the', 795), ('unt', 653), ('nto', 598), ('her', 590), ('aid', 533), ('sai', 488), ('ing', 48

```

1 print(f'Longitud del texto 1: {len(texto_1)}')
2 print(f'Longitud de bigramas del texto 1: {len(bigramas_texto_1)}')
3 print(f'Longitud de trigramas de caracteres del texto 1: {len(ngramas_caracteres_texto_1)}')

```

Longitud del texto 1: 110482
 Longitud de bigramas del texto 1: 110481
 Longitud de trigramas de caracteres del texto 1: 425131

```

1 print(f'Longitud del texto 2: {len(texto_2)}')
2 print(f'Longitud de bigramas del texto 2: {len(bigramas_texto_2)}')
3 print(f'Longitud de trigramas de caracteres del texto 2: {len(ngramas_caracteres_texto_2)}')

```

Longitud del texto 2: 53987
 Longitud de bigramas del texto 2: 53986
 Longitud de trigramas de caracteres del texto 2: 228150

```
1 print(f'Longitud del texto 3: {len(texto_3)}')
2 print(f'Longitud de bigramas del texto 3: {len(bigramas_texto_3)}')
3 print(f'Longitud de trigramas de caracteres del texto 3: {len(ngramas_caracteres_texto_3)}')
```

Longitud del texto 3: 18335
Longitud de bigramas del texto 3: 18334
Longitud de trigramas de caracteres del texto 3: 56086

▼ Similitud coseno

```
1 def crear_tabla_basica(nombre_texto, bigramas, trigramas):
2     df = pd.DataFrame()
3     bigramas_trigramas = bigramas + trigramas
4     df['texto'] = [nombre_texto for i in bigramas_trigramas]
5     df['bigramas y trigramas'] = [ngrama for ngrama in bigramas_trigramas]
6     return df
```

```
1 tabla_1 = crear_tabla_basica('texto_1', bigramas_texto_1,
2                               ngramas_caracteres_texto_1)
3 tabla_1.head(10)
```

	texto	bigramas y trigramas
0	texto_1	(moby, dick)
1	texto_1	(dick, herman)
2	texto_1	(herman, melville)
3	texto_1	(melville, etymology)
4	texto_1	(etymology, supplied)
5	texto_1	(supplied, late)
6	texto_1	(late, consumptive)
7	texto_1	(consumptive, usher)
8	texto_1	(usher, grammar)
9	texto_1	(grammar, school)

```
1 tabla_2 = crear_tabla_basica('texto_2', bigramas_texto_2,
2                               ngramas_caracteres_texto_2)
3 tabla_2.head(10)
```

	texto	bigramas y trigramas
0	texto_2	(sense, sensibility)
1	texto_2	(sensibility, jane)
2	texto_2	(jane, austen)
3	texto_2	(austen, chapter)

```
1 tabla_3 = crear_tabla_basica('texto_3', bigramas_texto_3,
2                               ngramas_caracteres_texto_3)
3 tabla_3.head(10)
```

	texto	bigramas y trigramas
0	texto_3	(beginning, god)
1	texto_3	(god, created)
2	texto_3	(created, heaven)
3	texto_3	(heaven, earth)
4	texto_3	(earth, earth)
5	texto_3	(earth, without)
6	texto_3	(without, form)
7	texto_3	(form, void)
8	texto_3	(void, darkness)
9	texto_3	(darkness, upon)

```
1 tabla_final = pd.concat([tabla_1, tabla_2, tabla_3], ignore_index=True)
2 tabla_final.tail(10)
```

	texto	bigramas y trigramas
892158	texto_3	lme
892159	texto_3	med
892160	texto_3	put
892161	texto_3	cof
892162	texto_3	off
892163	texto_3	ffi
892164	texto_3	fin
892165	texto_3	egy
892166	texto_3	gyp
892167	texto_3	ypt

```
1 tabla_pivoteada = tabla_final.groupby(['texto', 'bigramas y trigramas'])['bigramas y trigramas']\
2                               .agg(['count']).reset_index()\
3                               .pivot(index='bigramas y trigramas', columns='texto', values='count')
4
```

```

4
5 tabla_pivoteada.columns.name = None
6 tabla_pivoteada = tabla_pivoteada.fillna(0)
7 tabla_pivoteada[153477:153497]

```

	texto_1	texto_2	texto_3
bigramas y trigramas			
(zone, zone)	1.0	0.0	0.0
(zoned, letters)	1.0	0.0	0.0
(zoned, quest)	1.0	0.0	0.0
(zones, accounts)	1.0	0.0	0.0
(zones, surface)	1.0	0.0	0.0
(zones, world)	1.0	0.0	0.0
(zoology, anatomy)	1.0	0.0	0.0
(zoology, much)	1.0	0.0	0.0
(zoroaster, died)	1.0	0.0	0.0
(zuzims, ham)	0.0	0.0	1.0
	8.0	2.0	0.0
l	0.0	1.0	0.0
s	2.0	0.0	0.0
t	8.0	0.0	0.0
nd	1.0	0.0	0.0
st	4.0	0.0	0.0
th	10.0	0.0	0.0
—	22.0	0.0	0.0
aac	1.0	0.0	78.0
aah	0.0	0.0	1.0

```

1 from scipy.spatial.distance import cosine
2
3 def similitud_coseno(a,b):
4     distancia = cosine(a,b)
5     return 1-distancia
6
7 tabla_pivoteada.corr(method=similitud_coseno)

```

	texto_1	texto_2	texto_3
texto_1	1.000000	0.856819	0.644977
texto_2	0.856819	1.000000	0.600518
texto_3	0.644977	0.600518	1.000000

```

1 !wget -nc https://raw.githubusercontent.com/brpy/colab-pdf/master/colab_pdf.py

```

```
2 from colab_pdf import colab_pdf
3 colab_pdf('Tarea 10 - NLP.ipynb')
```

File 'colab_pdf.py' already there; not retrieving.

WARNING: apt does not have a stable CLI interface. Use with caution in scripts.

WARNING: apt does not have a stable CLI interface. Use with caution in scripts.

```
[NbConvertApp] Converting notebook /content/drive/My Drive/Colab Notebooks/Tarea 10 - NLP.ipynb
[NbConvertApp] Writing 50717 bytes to ./notebook.tex
[NbConvertApp] Building PDF
[NbConvertApp] Running xelatex 3 times: [u'xelatex', u'./notebook.tex', '-quiet']
[NbConvertApp] Running bibtex 1 time: [u'bibtex', u'./notebook']
[NbConvertApp] WARNING | bibtex had problems, most likely because there were no citations
[NbConvertApp] PDF successfully created
[NbConvertApp] Writing 43911 bytes to /content/drive/My Drive/Tarea 10 - NLP.pdf
'File ready to be Downloaded and Saved to Drive'
```