

Tarea 8

NLP

Descripción: Dados 3 textos construir la matriz término-documento y determinar la similitud coseno.

Alumno: Miguel Angel Soto Hernandez

Importaciones necesarias

```
# tratamiento de datos
import numpy as np
import pandas as pd
import string
import re

# preprocesado
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords

[ ]> [nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

Datos

```
# lectura de datos
url = 'https://raw.githubusercontent.com/JoaquinAmatRodrigo/Estadistica-con-R/master/datos/'
tweets_elon = pd.read_csv(url + "datos_tweets_@elonmusk.csv")
tweets_edlee = pd.read_csv(url + "datos_tweets_@mayoredlee.csv")
tweets_bgates = pd.read_csv(url + "datos_tweets_@BillGates.csv")

print('Número de tweets @BillGates: ' + str(tweets_bgates.shape[0]))
print('Número de tweets @mayoredlee: ' + str(tweets_edlee.shape[0]))
print('Número de tweets @elonmusk: ' + str(tweets_elon.shape[0]))

Número de tweets @BillGates: 2087
Número de tweets @mayoredlee: 2447
Número de tweets @elonmusk: 2678
```

tweets_elon.head()

	screen_name	user_id	created_at	status_id	text	retweet_count	favorite_count	is_quote_status	qu
0	elonmusk	44196397	2017-11-09T17:28:57Z	9.286758e+17	"If one day, my words are against science, cho...	49919	104722	False	
1	elonmusk	44196397	2017-11-09T17:12:46Z	9.286717e+17	I placed the flowers\n\nThree broken ribs\nA p...	5940	33725	False	
2	elonmusk	44196397	2017-11-08T18:55:13Z	9.283351e+17	Atatürk Anıtkabir https://t.co/al3wt0njr6	34752	104823	False	
3	elonmusk	44196397	2017-11-07T19:48:45Z	9.279862e+17	@Bob_Richards One rocket, slightly toasted	415	7247	False	
4	elonmusk	44196397	2017-10-28T21:36:18Z	9.243894e+17	@uncover007 500 ft so far. Should be 2 miles 1...	207	2128	False	

```
# se unen los dos dataframes en uno solo
tweets = pd.concat([tweets_elon, tweets_edlee, tweets_bgates], ignore_index=True)

# se seleccionan y renombran las columnas de interés
tweets = tweets[['screen_name', 'created_at', 'status_id', 'text']]
tweets.columns = ['autor', 'fecha', 'id', 'texto']
```

```
# parseo de fechas
tweets['fecha'] = pd.to_datetime(tweets['fecha'])
tweets.head()
```

	autor	fecha	id	texto
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	"If one day, my words are against science, cho...
1	elonmusk	2017-11-09 17:12:46+00:00	9.286717e+17	I placed the flowers\n\nThree broken ribs\nA p...
2	elonmusk	2017-11-08 18:55:13+00:00	9.283351e+17	Atatürk Anıtkabir https://t.co/al3wt0njr6
3	elonmusk	2017-11-07 19:48:45+00:00	9.279862e+17	@Bob_Richards One rocket, slightly toasted
4	elonmusk	2017-10-28 21:36:18+00:00	9.243894e+17	@uncover007 500 ft so far. Should be 2 miles l...

Limpeza y tokenizacion

```
def limpiar_tokenizar(texto):
    '''
    Esta función limpia y tokeniza el texto en palabras individuales.
    El orden en el que se va limpiando el texto no es arbitrario.
    El listado de signos de puntuación se ha obtenido de: print(string.punctuation)
    y re.escape(string.punctuation)
    '''

    # se convierte todo el texto a minúsculas
    nuevo_texto = texto.lower()

    # eliminación de páginas web (palabras que empiezan por "http")
    nuevo_texto = re.sub('http\S+', ' ', nuevo_texto)

    # eliminación de signos de puntuación
    regex = '[\!\@\#\$\%\&\'\(\)\*\+,\|\-\.\/\:\;\<\>=\>\|\?\\@\\[\]\^\_\`\\\{\|\}\|\~]'
    nuevo_texto = re.sub(regex, ' ', nuevo_texto)

    # eliminación de números
    nuevo_texto = re.sub("\d+", ' ', nuevo_texto)

    # eliminación de espacios en blanco múltiples
    nuevo_texto = re.sub("\s+", ' ', nuevo_texto)

    # tokenización por palabras individuales
    nuevo_texto = nuevo_texto.split(sep = ' ')

    # eliminación de tokens con una longitud < 2
    nuevo_texto = [token for token in nuevo_texto if len(token) > 1]

    return(nuevo_texto)
```

```
# prueba
test = '#HéroesSinCapa | El pasado 13 de marzo los oficiales Juárez, Morales y Castañeda adscritos a la #PBI, auxiliaron a los pacientes'
print(f'Ejemplo: {test}')
print(f'\nTexto limpio y por tokens: {limpiar_tokenizar(texto=test)}')
```

Ejemplo: #HéroesSinCapa | El pasado 13 de marzo los oficiales Juárez, Morales y Castañeda adscritos a la #PBI, auxiliaron a los pacientes

Texto limpio y por tokens: ['héroessincapa', 'el', 'pasado', 'de', 'marzo', 'los', 'oficiales', 'juárez', 'morales', 'castañeda', 'adscritos', 'a', 'la', '#pbi', 'auxiliaron', 'a', 'los', 'pacientes']

```
# se aplica la función de limpieza y tokenización a cada tweet
tweets['texto_tokenizado'] = tweets['texto'].apply(lambda x: limpiar_tokenizar(x))
tweets.head()
```

	autor	fecha	id	texto	texto_tokenizado
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	"If one day, my words are against science, cho...	[if, one, day, my, words, are, against, scienc...
1	elonmusk	2017-11-09 17:12:46+00:00	9.286717e+17	I placed the flowers\n\nThree broken ribs\nA p...	[placed, the, flowers, three, broken, ribs, pi...
2	elonmusk	2017-11-08 18:55:13+00:00	9.283351e+17	Atatürk Anıtkabir https://t.co/al3wt0njr6	[atatürk, anıtkabir]
3	elonmusk	2017-11-07 19:48:45+00:00	9.279862e+17	@Bob_Richards One rocket, slightly	[bob, richards, one, rocket, slightly,

```
# transformando los arreglos en un solo string por fila
```

```
# transformando los arreglos en un solo string por fila
# explode: transformar cada elemento de una lista en una fila,
# replicando los valores del índice
tweets_tidy = tweets.explode(column='texto_tokenizado')
tweets_tidy = tweets_tidy.drop(columns='texto')
tweets_tidy = tweets_tidy.rename(columns={'texto_tokenizado':'token'})
tweets_tidy.head()
```

	autor	fecha	id	token
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	if
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	one
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	day
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	my
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	words

```
# obtención de listado de stopwords del inglés
palabras_auxiliares = list(stopwords.words('english'))
```

```
# se añade la stopword: amp, ax, ex
palabras_auxiliares.extend(("amp", "xa", "xe"))
print(f'Palabras auxiliares: {palabras_auxiliares[:20]}')
```

Palabras auxiliares: ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd",

```
# limpiar palabras auxiliares
tweets_tidy = tweets_tidy[~(tweets_tidy['token'].isin(palabras_auxiliares))]
tweets_tidy.head()
```

	autor	fecha	id	token
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	one
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	day
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	words
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	science
0	elonmusk	2017-11-09 17:28:57+00:00	9.286758e+17	choose

▼ Frecuencia de palabras

```
# palabras totales por usuario
print('Palabras totales por usuario')
tweets_tidy.groupby(by='autor')['token'].count()
```

```
Palabras totales por usuario
autor
BillGates    19445
elonmusk     21719
mayoredlee   26676
Name: token, dtype: int64
```

```
# palabras unicas por usuario
print('Palabras distintas por usuario')
tweets_tidy.groupby(by='autor')['token'].nunique()
```

```
Palabras distintas por usuario
autor
BillGates    4718
elonmusk     6496
mayoredlee   5644
Name: token, dtype: int64
```

▼ Relacion coseno o correlacion

```
# pivotado de datos o matriz termino-documento
tweets_pivot = tweets_tidy.groupby(["autor","token"])["token"] \
    .agg(["count"]).reset_index() \
```

```
.pivot(index = "token" , columns="autor", values= "count")
tweets_pivot.columns.name = None
tweets_pivot = tweets_pivot.fillna(0)
tweets_pivot[1000:1030]
```

	BillGates	elonmusk	mayoredlee
token			
beliefs	0.0	0.0	3.0
believe	12.0	14.0	4.0
believed	1.0	1.0	0.0
believen	0.0	0.0	9.0
believing	0.0	3.0	0.0
bell	0.0	0.0	1.0
bellevue	0.0	0.0	1.0
bells	0.0	1.0	1.0
belltower	0.0	1.0	0.0
belongs	1.0	0.0	1.0
beloved	0.0	0.0	2.0
belovedrevol	0.0	2.0	0.0
benbenwilde	0.0	1.0	0.0
bencam	0.0	1.0	0.0
benchmark	0.0	1.0	0.0
bending	2.0	0.0	0.0
benedictevans	0.0	1.0	0.0
benefical	0.0	1.0	0.0
benefit	1.0	2.0	8.0
benefiting	1.0	0.0	0.0
benefits	13.0	0.0	10.0
benefitting	0.0	0.0	1.0
benefi...	0.0	0.0	1.0
benfeldman	0.0	1.0	0.0
benioff	0.0	0.0	6.0
benjamin	0.0	0.0	1.0
benjamincoop	0.0	1.0	0.0
benjaminrphoto	0.0	1.0	0.0
benmacy	0.0	1.0	0.0
benny	0.0	1.0	0.0

```
# relación coseno por el uso y frecuencia de palabras
from scipy.spatial.distance import cosine

def similitud_coseno(a,b):
    distancia = cosine(a,b)
    return 1-distancia

tweets_pivot.corr(method=similitud_coseno)
```

	BillGates	elonmusk	mayoredlee
BillGates	1.000000	0.413110	0.279001
elonmusk	0.413110	1.000000	0.197927
mayoredlee	0.279001	0.197927	1.000000

```
# número de palabras comunes
palabras_elon = set(tweets_tidy[tweets_tidy.autor == 'elonmusk']['token'])
palabras_bill = set(tweets_tidy[tweets_tidy.autor == 'BillGates']['token'])
```

```
palabras_edlee = set(tweets_tidy[tweets_tidy.autor == 'mayoredlee']['token'])
```

```
print(f"Palabras comunes entre Elon Musk y Ed Lee:\n      {len(palabras_elon.intersection(palabras_edlee))}")  
print(f"Palabras comunes entre Elon Musk y Bill Gates:\n      {len(palabras_elon.intersection(palabras_bill))}")  
print(f"Palabras comunes entre Bill Gates y Ed Lee:\n      {len(palabras_bill.intersection(palabras_edlee))}")
```

```
Palabras comunes entre Elon Musk y Ed Lee:      1760  
Palabras comunes entre Elon Musk y Bill Gates:   1758  
Palabras comunes entre Bill Gates y Ed Lee:     1717
```

✓ 0s completed at 8:05 PM

