

```
# BETZABETH MADELINE MUÑOZ VILLEGAS
```

```
# Tratamiento de datos
```

```
# =====
```

```
import numpy as np
import pandas as pd
import string
import re
```

```
# Gráficos
```

```
# =====
```

```
import matplotlib.pyplot as plt
from matplotlib import style
import seaborn as sns
#style.use('ggplot') or plt.style.use('ggplot')
```

```
# Preprocesado y modelado
```

```
# =====
```

```
from sklearn import svm
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import confusion_matrix
from sklearn.feature_extraction.text import TfidfVectorizer
import nltk
#nltk.download('stopwords')
from nltk.corpus import stopwords
```

```
# Configuración warnings
```

```
# =====
```

```
import warnings
warnings.filterwarnings('ignore')
```

```
from google.colab import drive
drive.mount('/content/drive')
```

```
# Lectura de datos --- BETZABETH MADELINE MUÑOZ VILLEGAS
```

```
# =====
```

```
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
```

```
# 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(3)
```

	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

Distribución temporal de los tweets --- BETZABETH MADELINE MUÑOZ VILLEGAS

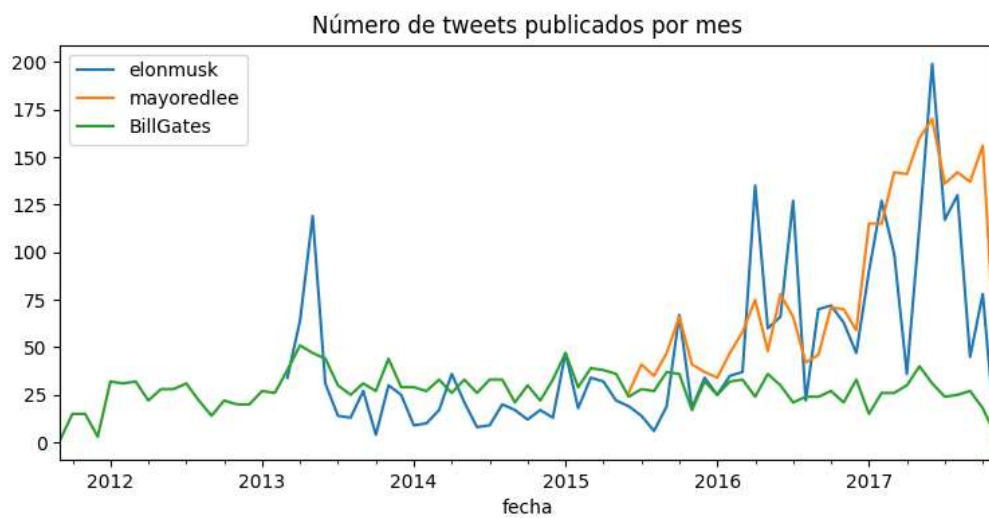
=====

```
fig, ax = plt.subplots(figsize=(9,4))
```

```
for autor in tweets.autor.unique():
    df_temp = tweets[tweets['autor'] == autor].copy()
    df_temp['fecha'] = pd.to_datetime(df_temp['fecha']).dt.strftime('%Y-%m')
    df_temp = df_temp.groupby(df_temp['fecha']).size()
    df_temp.plot(label=autor, ax=ax)
```

```
ax.set_title('Número de tweets publicados por mes')
```

```
ax.legend();
```



```
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)
```

```
test = "Esto es 1 ejemplo de l'limpieza de6 TEXTO https://t.co/rnHPgyhx4Z @cienciadedatos #textmining"
```

```
print(test)
print(limpiar_tokenizar(texto=test))
```

Esto es 1 ejemplo de l'limpieza de6 TEXTO <https://t.co/rnHPgyhx4Z> @cienciadedatos #textmining
 ['esto', 'es', 'ejemplo', 'de', 'limpieza', 'de', 'texto', 'cienciadedatos', 'textmining']

Palabras más utilizadas

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

	texto	texto_tokenizado
0	"If one day, my words are against science, cho...	[if, one, day, my, words, are, against, scienc...
1	I placed the flowers\n\nThree broken ribs\nA p...	[placed, the, flowers, three, broken, ribs, pi...
2	Atatürk Anıtkabir https://t.co/al3wt0njr6	[atatürk, anıtkabir]
3	@Bob_Richards One rocket, slightly toasted	[bob, richards, one, rocket, slightly, toasted]
4	@uncover007 500 ft so far. Should be 2 miles	[uncover, ft, so, far, should, be, miles, long]

```
# Unnest de la columna texto_tokenizado
# =====
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(3)
```

	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		

```
# Top 5 palabras más utilizadas por cada autor
# =====
tweets_tidy.groupby(['autor', 'token'])['token'] \
    .count() \
    .reset_index(name='count') \
    .groupby('autor') \
    .apply(lambda x: x.sort_values('count', ascending=False).head(5))
```

		autor	token	count		
		autor				
	BillGates	4195	BillGates	the	1178	
		4271	BillGates	to	1115	
		2930	BillGates	of	669	
		2084	BillGates	in	590	

Frecuencia de palabras

```

# Palabras totales utilizadas por cada autor
# =====
print('-----')
print('Palabras totales por autor')
print('-----')
tweets_tidy.groupby(by='autor')['token'].count()

-----
Palabras totales por autor
-----
autor
BillGates      31500
elonmusk       33609
mayoredlee     41878
Name: token, dtype: int64

```

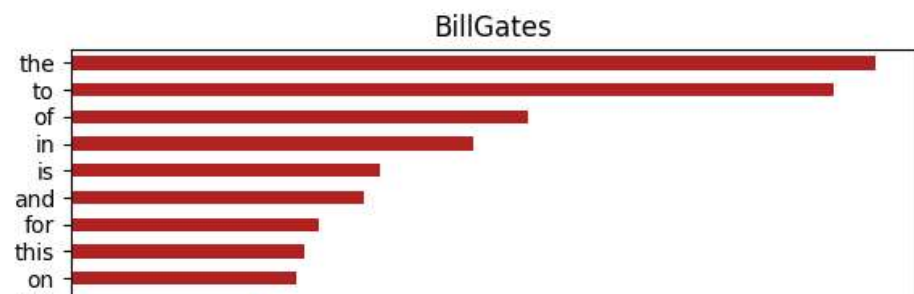
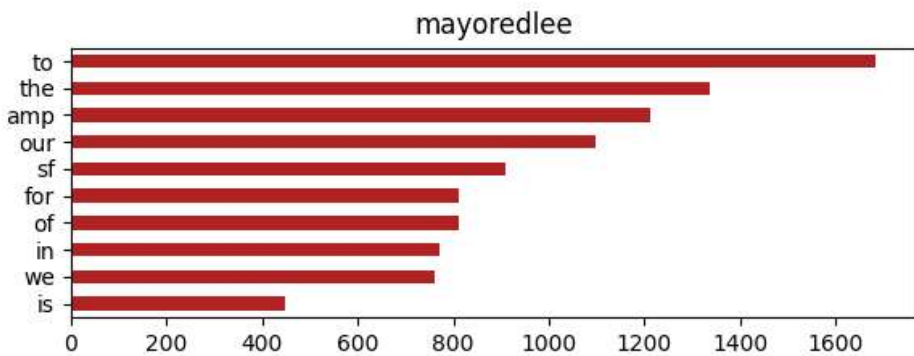
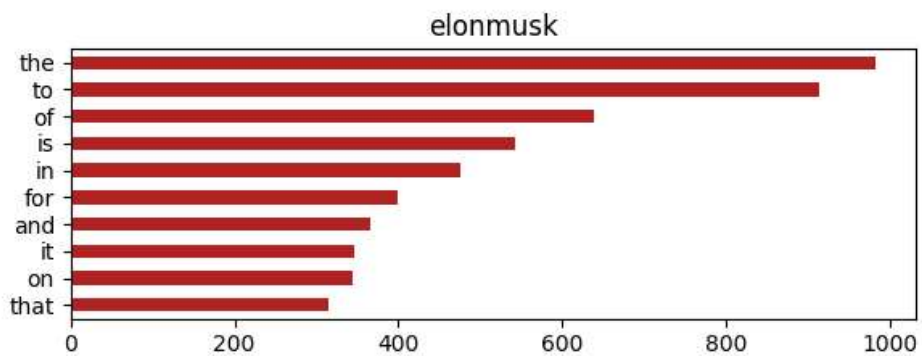
Terminos mas frecuentes en todos los usuarios

```

# Top 10 palabras por autor (sin stopwords)
# =====
fig, axs = plt.subplots(nrows=3, ncols=1, figsize=(6, 7))
for i, autor in enumerate(tweets_tidy.autor.unique()):
    df_temp = tweets_tidy[tweets_tidy.autor == autor]
    counts = df_temp['token'].value_counts(ascending=False).head(10)
    counts.plot(kind='barh', color='firebrick', ax=axs[i])
    axs[i].invert_yaxis()
    axs[i].set_title(autor)

fig.tight_layout()

```



✓ 1 s se ejecutó 21:03

● ✕