

```

1 # Descarga de archivos de ProfNER
2 !wget https://zenodo.org/record/4563995/files/profner.zip?download=1
3 # Si el unzip no funciona, revisar cual es el nombre de descarga del archivo
4 !unzip profner.zip?download=1

```

Funciones de carga y preparación de datos

+ Code

+ Text

```

1 import sys
2 import os
3 import pandas as pd
4 def get_tweet_content(list_paths):
5     """
6     Función para guardar en un diccionario el contenido de archivos txt
7     que se introduce en su entrada.
8     Devuelve un diccionario en el que las claves son el id del tweet, y
9     el valor el texto del tweet.
10    """
11    output_dict = dict()
12    for i in list_paths:
13        tweet_id = i.split("/")[-1].split(".txt")[0]
14        with open(i) as f:
15            output_dict[int(tweet_id)] = f.read()
16
17    return output_dict
18
19 def get_profner_data(profner_path_data):
20     # Obtenemos el path a los txt de los tweets.
21     path_to_txt = profner_path_data+"subtask-1/train-valid-txt-files/"
22     tweets_train_files = [path_to_txt+"train/"+i for i in os.listdir(path_to_txt+"train/")]
23     tweets_valid_files = [path_to_txt+"valid/"+i for i in os.listdir(path_to_txt+"valid/")]
24     # Obtenemos diccionarios en los que el key es el tweet_id y el value el texto del tweet.
25     train_txt_content = get_tweet_content(tweets_train_files)
26     valid_txt_content = get_tweet_content(tweets_valid_files)
27
28     # Cargamos dos dataframes con los tweet_id y la categoría de los tweets
29     path_to_labeled = profner_path_data+"subtask-1/"
30     train_tweets = pd.read_csv(path_to_labeled+"train.tsv", sep="\t")
31     valid_tweets = pd.read_csv(path_to_labeled+"valid.tsv", sep="\t")
32
33     # Introducimos a los df el campo de texto mapeando los diccionarios con tweet_id
34     train_tweets["tweet_text"] = train_tweets['tweet_id'].map(train_txt_content)
35     train_tweets["set"] = "train"
36     valid_tweets["tweet_text"] = valid_tweets['tweet_id'].map(valid_txt_content)
37     valid_tweets["set"] = "valid"
38
39     # Concatenamos el resultado
40     output_df = pd.concat([train_tweets, valid_tweets], axis=0)
41     # Eliminamos retorno de carro
42     output_df["tweet_text"] = output_df.tweet_text.apply(lambda x: x.replace('\n', ' '))
43     return output_df[["tweet_id", "tweet_text", "label", "set"]].reset_index(drop=True)

```

▼ Ejercicio

En este ejercicio se trabajará con un conjunto de datos reales publicados para la shared-task [ProfNER](#), celebrada en el año 2021.

Específicamente, se utilizarán los datos textuales de la subtask 1, centrada en la clasificación de textos. Este conjunto de datos son tweets en español que tienen asignada una etiqueta numérica, que representa la presencia (valor 1) o no (valor 0) de menciones de profesiones en el tweet. Por si fuera de tu interés, el proceso de obtención, selección y anotación de datos está descrita en [este enlace](#).

Para el ejercicio debéis entrenar diferentes modelos de clasificación que permitan clasificar correctamente los tweets. Para ello será necesario crear y utilizar funciones de preprocesado de datos similares a las vistas en clase, aplicar estrategias de vectorización de textos como TF-IDF o embeddings, y entrenar/evaluar modelos de clasificación. Para que os sirva de orientación, los criterios de evaluación del ejercicio serán los siguientes:

- **Análisis exploratorio, pre-procesado y normalización de los datos (30%):**
 - El ejercicio deberá contener un análisis exploratorio de los datos como número de documentos, gráficas de distribución de longitudes y/o wordclouds, entre otros análisis que se os pudieran ocurrir. Vuestros ejercicios deberán incorporar al menos los análisis exploratorios vistos en clase.
 - También tendréis que tener funciones para normalizar textos que permitan eliminar palabras vacías, quitar símbolos de puntuación y lematizar o hacer stemming.
- **Vectorización de textos (40%)**

En clase hemos visto diferentes estrategias de vectorización como TF-IDF y Word Embeddings. También hemos visto como incorporar características adicionales utilizando el sentimiento de los documentos. Para este ejercicio sois libres de utilizar la estrategia de vectorización que queráis, pero:

- Si decidís utilizar TF-IDF será necesarios que incorporéis a modelo características adicionales de sentimiento utilizando recursos adicionales (como por ejemplo la librería TextBlob).
 - Si optáis por representar el texto mediante embeddings, dado que en clase no se profundizó sobre el tema no será necesario incorporar esas características adicionales. Si decidís esta segunda opción, podéis utilizar los embeddings en español que vimos en clase
- **Entrenamiento y validación del sistema (30%)**
 - En el proceso de entrenamiento del modelo tendréis que testear al menos 3 modelos de clasificación. El procedimiento debe ser similar al visto en clase, en el que primero estimábamos el rendimiento de varios algoritmos de forma general, para posteriormente seleccionar el mejor para ajustar los hiperparámetros.

▼ 0. Imports

```

1 # Instalamos nltk
2 !pip install nltk
3 !pip install contractions
4
5 # Importamos
6 import nltk
7 # Complementos de la librería necesarios para su funcionamiento.
8 # Todas las opciones aquí https://www.nltk.org/nltk_data/
9 nltk.download('punkt')
10 nltk.download('wordnet')
11 nltk.download('averaged_perceptron_tagger')
12 nltk.download('tagsets')
13 nltk.download('maxent_ne_chunker')
14 nltk.download('words')
15 nltk.download('stopwords')
16 !pip install emosent-py
17 !pip install emoji_extractor
18 !pip install emoji
19 !pip install vaderSentiment
20 from textblob import TextBlob
21 !wget https://www.clarin.si/repository/xmlui/handle/11356/1048/allzip
22 !unzip allzip
23 import nltk
24 nltk.download('punkt')
25 nltk.download('averaged_perceptron_tagger')
26 nltk.download('opinion_lexicon')
27 nltk.download('subjectivity')
28 nltk.download('vader_lexicon')
29 nltk.download('wordnet')
30 # Instalamos textacy
31 !pip install textacy
32 # Instalamos spacy y uno de sus modelos
33 !pip install spacy = 3.2.1
34 # Descargamos modelos pre-entrenados de spacy.
35 !python -m spacy download en_core_web_md
36
37 import matplotlib.pyplot as plt
38 import sys
39 import numpy as np
40 import pandas as pd
41 import seaborn as sns
42 import en_core_web_sm
43 import contractions
44 from nltk.tokenize import TweetTokenizer
45 from nltk.corpus import stopwords
46 from sklearn.feature_extraction.text import TfidfVectorizer
47 from sklearn.linear_model import LogisticRegression
48 from sklearn.model_selection import train_test_split
49 from sklearn.metrics import f1_score, confusion_matrix
50 from sklearn.metrics import classification_report
51
52 from tqdm.notebook import tqdm
53 tqdm.pandas()

1 from nltk import tokenize

```

▼ 1. Obtención del corpus

Para la obtención de los datos teneis disponible la función `get_profner_data()` . Esta función prepara los datos del ejercicio en formato Pandas dataframe para que podais realizarlo.

```
1 profner_path = "./profner/"
2 datos_ejercicio = get_profner_data(profner_path)

1 datos_ejercicio.head(10)
```

	tweet_id	tweet_text	label	set
0	1256007275807997953	CHINA: libera una pandemia EE.UU: libera OVNIS...	0	train
1	1257909427677601792	San Francisco (EEUU) hace 100 años tras vivir ...	0	train
2	1281580572415066112	Porfi poneos la mascarilla o tendremos 28 nuev...	0	train
3	1271348112733470720	El nuevo „covid normas y reglas recibimiento“ ...	0	train
4	1270421287148695556	Si el confinamiento ha dejado algo tocada tu e...	0	train

```
1 datos_ejercicio['set'].value_counts()

train    6000
valid    2000
Name: set, dtype: int64
```

TRANSLATE BOOONA

```
1 !pip install -U deep-translator

1 from deep_translator import (GoogleTranslator,
2                               MicrosoftTranslator,
3                               PonsTranslator,
4                               LingueeTranslator,
5                               MyMemoryTranslator,
6                               YandexTranslator,
7                               PapagoTranslator,
8                               DeeplTranslator,
9                               QcriTranslator,
10                              single_detection,
11                              batch_detection)

1 tweet_cast = datos_ejercicio['tweet_text']

1 tweet_cast = pd.Series(tweet_cast)

1 tweet_cast = "\n\n\n\n\n\n\n".join(tweet_cast.astype(str).to_list())

1 tweet_cast

'CHINA: libera una pandemia EE.UU: libera OVNIS ARGENTINA: libera
presos Jajajajajajaja\n\n\n\n\n\n\nSan Francisco (EEUU) hace 100
años tras vivir la pandemia de 1918. https://t.co/veOT7nkdi8\n\n
\n\n\n\n\n\nPorfi poneos la mascarilla o tendremos 28 nuevas cancio

1 chunks = [tweet_cast[x:x+1755] for x in range(0, len(tweet_cast), 1755)]

1 translated = GoogleTranslator('es', 'en').translate_batch(chunks)

1 translated = str(translated)

1 translated = translated.split("\n\n\n\n\n\n\n\n\n\n\n\n\n\n\n\n")
```

1 translated

['CHINA: releases a pandemic USA: releases UFOs ARGENTINA: releases prisoners Hahahahahahaha',
 'San Francisco (USA) 100 years ago after experiencing the 1918 pandemic. <https://t.co/ve077nkd18>',
 'Please put on your mask or we will have 28 new songs from artists collaborating with Rozalén.',
 'The new „covid norms and reception rules“ in the Spanish retail trade. Walking from @mkdirecto one more day the streets and shops of Madrid to know and show the new "normal consumption": <https://t.co/EPrZ9Ncb5m>',
 'If confinement has left your mental stability somewhat affected, you can always console yourself by thinking that Miguel Bosé has been much worse off.',
 'What has never been seen, here in Spain, from what we can see, the dead resurrect us. I see less and less. I do not say it with irony, I say it with a lot of indignation, I qualify it / 27,125 deaths in Spain: the coronavirus, community by community <https://t.co/PBD6ktbWNV> Sent from @updayESP',
 '#Coronavirus: A new study finds that both symptomatic and asymptomatic Covid-19 patients have the ability to contaminate their surroundings. I mean, everything they touch <https://t.co/FELV3KPLPH>',
 'This will be an archaeological vestige of the ball the day after they cut the tape',
 'Delivery men having to work in the middle of a pandemic that has us all locked up and they get a €1,500 fine for parking on a yellow line, you have to be a shitty pig, without more.',
 '@carlesenric @salvadorilla De-escalation is essential, with health criteria, not political ones, with managers with at least a CV as in Italy, the population cannot be put at risk, seeking political support... we are facing a pandemic',
 'The recovery of employment prior to', 'COVID-19 did not stop the increase in severe poverty among the severely excluded population. The crisis and its push for unemployment will continue to plunge many families into severe poverty.
 #ElPoderDeCadaPersona <https://t.co/Yr9YeS0zZe>',
 'Coronavirus in Yucatan: 109 new infections and 16 deaths <https://t.co/AwkspwljLI> <https://t.co/3V69toRnhu>',
 'A new dawn Good morning! A tulip... ~haiku~ #LYF15 #VAFIorDePiel #YoMeQuedoEnCasa <https://t.co/L8nbM4cN20>',
 'The reasons why Madrid does not impose the mandatory use of a mask <https://t.co/wCTekphfDV>',
 'the same their Spanish antibodies are not so Spanish',
 'TVE Exhumation of Franco: -22 cameras -3 mobile units -50 people -Live broadcast Funeral of 45,000 Spaniards killed by coronavirus: -ZERO deployment -NO broadcast The Government wants you to remember that it unearthed a dead person, but not that it buried 45,000 .',
 'It is good for us fexs that it is mandatory to wear a mask',
 'The same effort that you make, that you don\\'t even open the door of your official car and collect your full salary, not like the millions of Spaniards affected by the ERTE or the self-employed that we pay without being paid, you are not and never were working class",
 'Context: - Police station with boxes of masks at the door - Politician and press come to record the event - In the first shot they realize that he is not wearing a mask - They repeat the shot - Boxes are put in the vehicle and they leave 📷 In whose hands we\\'re? <https://t.co/0cF8xZPpos>',
 'That the director of the National Epidemiology Center says "since February 1, that is, since we have the epidemic in our country" do we take it as fake? lapse? sologr\\', \\'ipism? It\\'s just that I ended up messing around. <https://t.co/ewd2Czc34g>',
 '@Awuamba I\\'m not saying it\\'s easy, but it wouldn\\'t be anything new for us, personally, I think confinement wouldn\\'t be as bad as I see people are doing. I also talk about what I know and my experience, there will be those who agree and those who do not.",
 '@ danigandara4 also sends us his message #StayHome <https://t.co/bKJXuwmnr>',
 'Several million euros spent on advertising to promote Interior tourism and Fernando Simón is caught in Portugal without a mask and without Social distance. Public regulations and hospitals are for the populace. #givingexample <https://t.co/7UDnWoEW7>',
 'Ambassador Buchan: The US Embassy in Madrid has projected the Spanish flag on its facade in solidarity with the Spanish people in the fight against #COVID19. In these dark times, the relationship between #USA and #Spain as partners, friends and allies continues to shineuses <https://t.co/beXpnEyNAJ>',
 'There is not a single word in this tweet that is not a lie.',
 '@Ernesto_Fado @Cascabullo48 @europapress Okay... I thought I was the only one who was caught with the little phrase... If it\\'s difficult... Don\\'t try it, I think there are people more prepared in other places, I don\\'t see that are capable of multitasking... Either pandemic or climate change, the most urgent thing is the pandemic...',
 '📌 April 7th #healthdiamond 📌 @WHO dedicates it to nurses, with the motto "Support nurses and midwives" highlights some goals #NursingNow 2020 📌 In a pandemic #Covid_19 your work is essential! 📌 #QuedateEnCasa #EsteVirusLoParamosUnidos @saludand <https://t.co/vRBpcUGXY9> <https://t.co/4HjvZ8hyGN>',
 'It is no longer that it is disloyal, it is that it shows an indi', 'Insulting intellectual agency: How can you put in the same tweet that Sánchez hides deaths and at the same time use the INE or the Carlos III Institute, which depend on the Sánchez government, to argue it? Spain deserves another right",
 '@ctabuyoiOS @Xixonto1 @Kusumoto_Ine @CRossanacoll @mavidonate @cfranganillo @telediarario_tve And how do you explain that in China, without seeing the pandemic coming, fewer people die than in other countries?',
 '► @IdiazAyuso in @elprogramadear: "We must go against this epidemic, prevent this from happening again. We must have a clear

▼ 2. Análisis exploratorio

EDA

Número de documentos y columnas

```
1 # Transformar la columna "text" a una lista
2 texto_noticias = pd.DataFrame(translated)
3 print(type(texto_noticias))
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
1 texto_noticias.rename(columns={0: 'texto_eng'},
2                       inplace=True)
```

```
1 print("El corpus texto_noticias contiene un total de {} documentos".format(len(texto_noticias)))
2 print("El dataframe tiene {} columnas".format(texto_noticias.shape))
```

```
El corpus texto_noticias contiene un total de 7974 documentos
El dataframe tiene (7974, 1) columnas
```

```

1 label1 = datos_ejercicio['label']

1 set1 = datos_ejercicio['set']

1 label1 = label1[0:7974]

1 set1 = set1[0:7974]

1 texto_noticias['label'] = label1

1 texto_noticias['set'] = set1

1 texto_noticias

```

	texto_eng	label	set
0	[CHINA: releases a pandemic USA: releases UFO...	0	train
1	San Francisco (USA) 100 years ago after experi...	0	train
2	Please put on your mask or we will have 28 new...	0	train
3	The new „covid norms and reception rules“ in t...	0	train
4	If confinement has left your mental stability ...	0	train
...
7969	#MostRead The new coronavirus infected multi...	0	valid
7970	Ayuso awards a contract of 30,000 euros by han...	0	valid
7971	months locked up You can travel to other provi...	1	valid
7972	The government prohibits layoffs by law during...	0	valid
7973	Psychological traits of why you don't wear a #...	0	valid

7974 rows × 3 columns

Podemos observar que al traducir el corpus ha pasado de 8000 filas a 7974. Se han perdido 26 filas test con algún valor 1 en label. A mi parecer es un error pero me ha facilitado la elaboración del ejercicio, por lo tanto, he decidido obviarlo ya que la pérdida de información ronda el 3.5%.

Número de documentos duplicados:

```

1 print("Existen {} noticias duplicadas".format(np.sum(texto_noticias.duplicated(subset=["texto_eng"]))))
2 # Quitaremos esos duplicados
3 texto_noticias = texto_noticias.drop_duplicates()
4 print("Despues de quitar duplicados tenemos un conjunto de {} noticias".format(texto_noticias.shape[0]))

Existen 0 noticias duplicadas
Despues de quitar duplicados tenemos un conjunto de 7974 noticias

1 print("Hay {} valores vacíos en las noticias y {} valores vacíos en las etiquetas en los datos".format(np.sum(texto_noticias.isnull()))

Hay 0 valores vacíos en las noticias y 0 valores vacíos en las etiquetas en los datos

```

Número de documentos por cada clase:

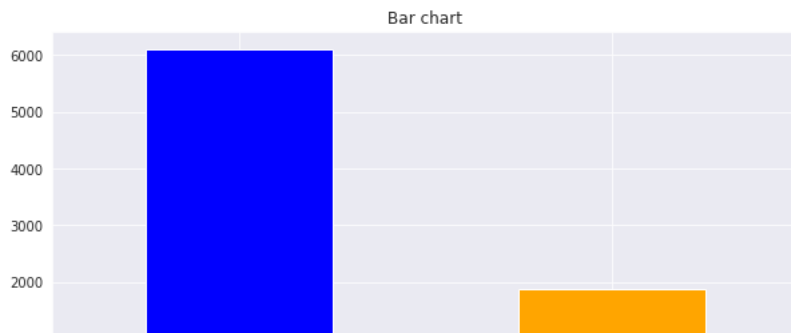
```

1 texto_noticias['label'].value_counts()

0    6109
1    1865
Name: label, dtype: int64

1 ax, fig = plt.subplots()
2 etiquetas = texto_noticias.label.value_counts()
3 etiquetas.plot(kind='bar', color = ["blue", "orange"])
4 plt.title('Bar chart')
5 plt.show()

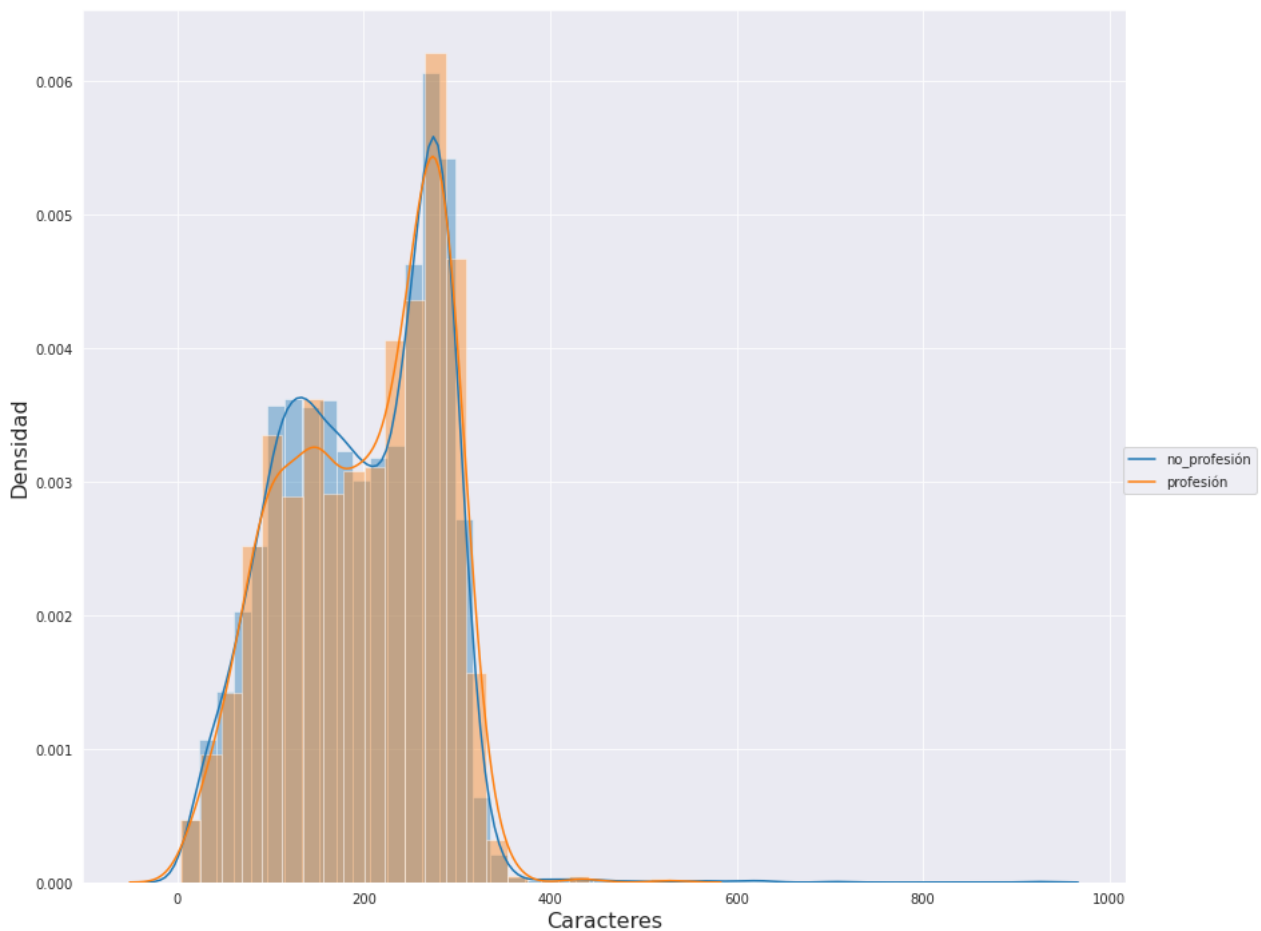
```



Distribución de la longitud de los tweets en caracteres:

```
0
1 texto_noticias['char_len'] = texto_noticias['texto_eng'].apply(lambda x: len(x))
```

```
1 # Importamos las librerías matplotlib y seaborn:
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 fig = plt.figure(figsize=(14,12))
6 sns.set_style("darkgrid")
7 # añadimos series para cada categoría (eligiendo la serie de char_len_ruido)
8 plt1 = sns.distplot(texto_noticias[texto_noticias.label==0].char_len, hist=True, label="no_profesión")
9 plt2 = sns.distplot(texto_noticias[texto_noticias.label==1].char_len, hist=True, label="profesión")
10 fig.legend(labels=['no_profesión', 'profesión'], loc = 5)
11
12
13 # Definimos el título de los ejes:
14 plt.xlabel('Caracteres', fontsize=16)
15 plt.ylabel('Densidad', fontsize=16)
16
17 # Finalmente mostramos el gráfico:
18 plt.show()
```



Podemos observar como la distribución de las longitud de los tweets es muy similar.

Wordcloud (antes de preprocesar)

[illegible]

Normalización del texto

7/14

```

1 # Eliminar espacios
2 texto_noticias["normaliza"] = texto_noticias["texto_eng"].progress_apply(lambda x: eliminar_espacios(x))
3
4 # To Lower
5 texto_noticias["normaliza"] = texto_noticias["normaliza"].progress_apply(lambda x: texto_to_lower(x))
6
7 # Quitar Contractions
8 texto_noticias["normaliza"] = texto_noticias["normaliza"].progress_apply(lambda x: replace_contraction(x))
9
10 # Tokenizar
11 texto_noticias["normaliza"] = texto_noticias["normaliza"].progress_apply(lambda x: tokenize(x))
12
13 # Quitar Stopwords
14 texto_noticias["normaliza"] = texto_noticias["normaliza"].progress_apply(lambda x: quitar_stopwords(x))
15
16 # Quitar puntuación
17 texto_noticias["normaliza"] = texto_noticias["normaliza"].progress_apply(lambda x: quitar_puntuacion(x))
18
19 # Mirar todo lo que tarda con lematización (mediante spacy)
20 texto_noticias["normaliza"] = texto_noticias["normaliza"].progress_apply(lambda x: lematizar(x))
21
22
23 texto_noticias['normaliza']

```

```

100% 7974/7974 [00:00<00:00, 51115.17it/s]
100% 7974/7974 [00:00<00:00, 66559.15it/s]
100% 7974/7974 [00:01<00:00, 8181.34it/s]
100% 7974/7974 [00:03<00:00, 1679.55it/s]
100% 7974/7974 [00:05<00:00, 1874.17it/s]
100% 7974/7974 [00:00<00:00, 19593.18it/s]
100% 7974/7974 [01:35<00:00, 134.95it/s]
0      [china, release, pandemic, usa, release, ufos,...
1      [san, francisco, usa, 100, year, ago, experien...
2      [please, put, mask, 28, new, song, artist, col...
3      [new, covid, norm, reception, rule, spanish, r...
4      [confinement, leave, mental, stability, somewh...
...
7969   [new, coronavirus, infect, multiply, almost, 2...
7970   [ayuso, awards, contract, euros, hand, make, v...
7971   [month, lock, travel, province, go, restaurant...
7972   [government, prohibit, layoff, law, coronaviru...
7973   [psychological, trait, wear]
Name: normaliza, Length: 7974, dtype: object

```

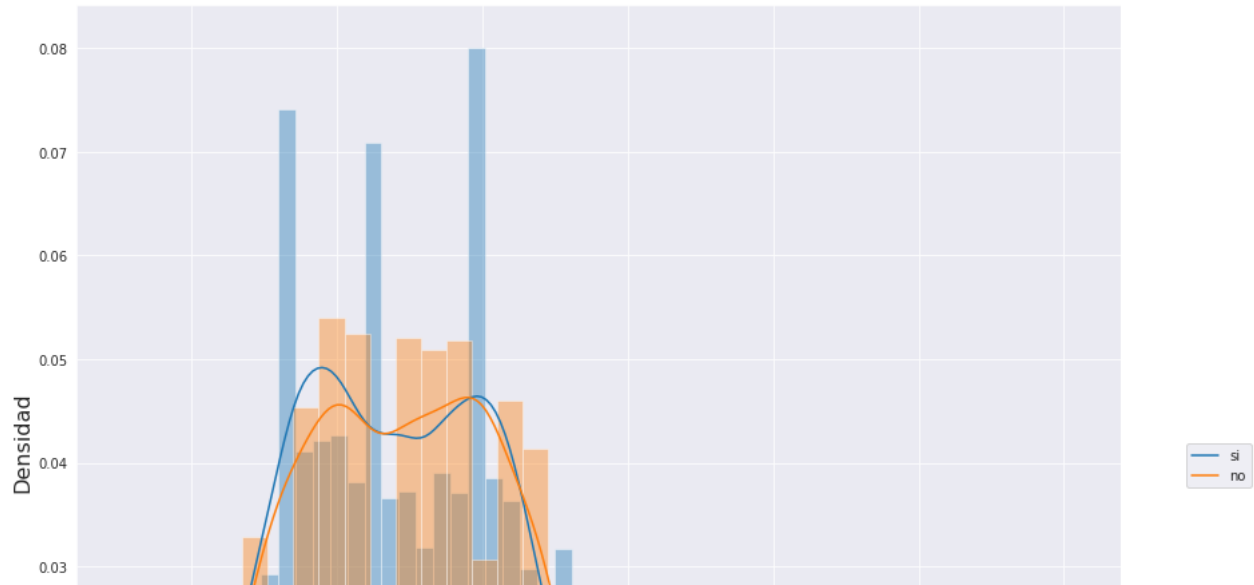
```
1 texto_noticias['clean_text'] = texto_noticias["normaliza"].progress_apply(lambda x: " ".join(x))
```

```
100% 7974/7974 [00:00<00:00, 89577.52it/s]
```

```

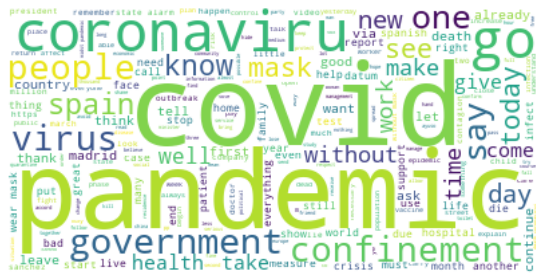
1 texto_noticias["token_len_limpio"] = texto_noticias["normaliza"].apply(lambda x: len(x))
2
3 fig = plt.figure(figsize=(14,12))
4 sns.set_style("darkgrid")
5 plt1 = sns.distplot(texto_noticias[texto_noticias.label == 0].token_len_limpio, hist=True, label="no")
6 plt2 = sns.distplot(texto_noticias[texto_noticias.label == 1].token_len_limpio, hist=True, label="si")
7 fig.legend(labels=['si', 'no'], loc=5)
8
9 # Definimos el título de los ejes:
10 plt.xlabel('Tokens', fontsize=16)
11 plt.ylabel('Densidad', fontsize=16)
12
13 plt.show()

```

Vemos como la distribución de la lonitud de los tweets a sido variada a causa del pre-procesamiento. Podemos observar picos en la densidad que antes no teníamos.

```
1 from wordcloud import WordCloud
2 # Une las frases
3 long_string = ', '.join(list(texto_noticias["clean_text"].values))
4 # Genera un objeto WordCloud
5 wordcloud = WordCloud(background_color="white", max_words=5000, contour_width=0, contour_color='steelblue')
6 # Genera el wordcloud
7 wordcloud.generate(long_string)
8 # Visualizalo en una imagen
9 wordcloud.to_image()
```



Una vez preprocesado el corpus podemos observar como han cambiado las palabras que más se usan por palabras con significados relevantes.

	texto_eng	label	set	char_len	normaliza	clean_text	token_len_limpio
0	[CHINA: releases a pandemic USA: releases UFO...	0	train	92	[china, release, pandemic, usa, release, ufos,...	china release pandemic usa release ufos argent...	10
1	San Francisco (USA) 100 years ago after experi...	0	train	95	[san, francisco, usa, 100, year, ago, experien...	san francisco usa 100 year ago experience 1918...	9
2	Please put on your mask or we will have 28 new...	0	train	93	[please, put, mask, 28, new, song, artist, col...	please put mask 28 new song artist collaborate...	9
3	The new „covid norms and reception rules“ in t...	0	train	210	[new, covid, norm, reception, rule, spanish, r...	new covid norm reception rule spanish retail t...	19
4	If confinement has left your mental stability ...	0	train	150	[confinement, leave, mental, stability, somewh...	confinement leave mental stability somewhat af...	13
...
7969	#MostRead The new coronavirus infected multi...	0	valid	136	[new, coronavirus, infect, multiply, almost, 2...	new coronavirus infect multiply almost 20 time...	11
7970	Ayuso awards a contract of 30,000 euros by han...	0	valid	137	[ayuso, awards, contract, euros, hand, make, v...	ayuso awards contract euros hand make videos c...	11
...

Visualización

```

1 def frecuencia_tokens(lista):
2     # Creamos diccionario vacío
3     frecuencia = {}
4     for item in lista:
5         if (item in frecuencia):
6             frecuencia[item] += 1
7         else:
8             frecuencia[item] = 1
9     return frecuencia
10
11 from nltk.corpus.reader.tagged import word_tokenize
12 lista_tokens = list()
13 for i in texto_noticias['clean_text']:
14     # Tokenizamos cada documento con word_tokenize()
15     tokens_document = word_tokenize(i)
16     # Añadimos esos tokens como nuevos elementos
17     # Si usamos append se crearía una lista de listas, de este modo añadimos los
18     lista_tokens.extend(tokens_document)
19
20

```

```

1 from collections import Counter

```

```

1 %%time
2 dict_freq = Counter(lista_tokens)
3 dict_freq["Road"]

CPU times: user 22.6 ms, sys: 0 ns, total: 22.6 ms
Wall time: 22.4 ms
0

```

```

1 # Ordenamos el diccionario por la frecuencia de sus palabras
2 dict_freq_order = sorted(dict_freq.items(), key=lambda x: x[1], reverse=True)
3 token_names = list()
4 token_freqs = list()
5 for i in dict_freq_order:
6     if i[1] > 350:
7         token_names.append(i[0])
8         token_freqs.append(i[1])

```

```

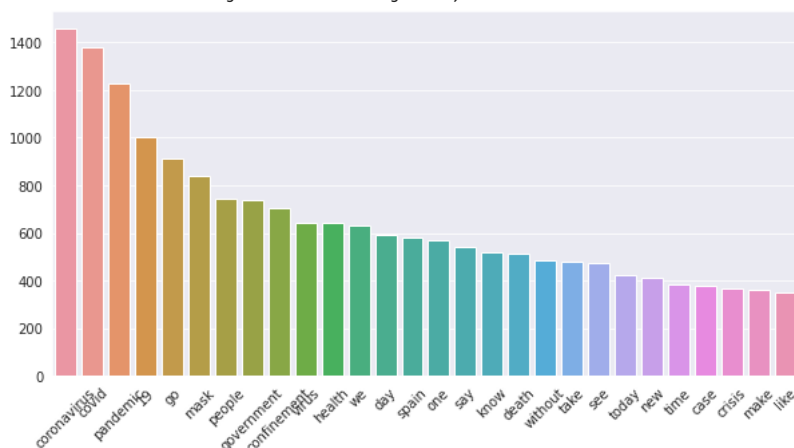
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 plt.rcParams['figure.figsize'] = [10, 5]
4 sns_g = sns.barplot(x=token_names, y=token_freqs)
5 plt.xticks(rotation=45)

```

```

(array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16,
        17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27]),
<a list of 28 Text major ticklabel objects>)

```



Extracción de sentimiento

```

1 from textblob import TextBlob

```

```

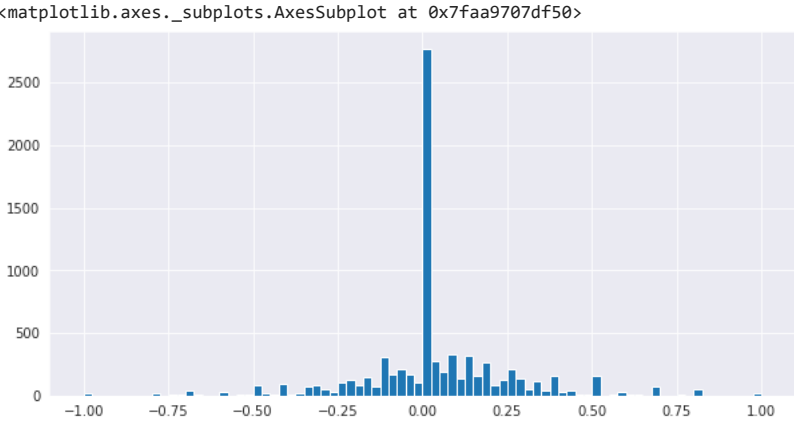
1 texto_noticias["sent_subjectivity_str"] = texto_noticias["clean_text"].progress_apply(lambda x: TextBlob(x).sentiment.subjectivity)
2 texto_noticias["sent_polarity_str"] = texto_noticias["clean_text"].progress_apply(lambda x: TextBlob(x).sentiment.polarity)

```

100%

7974/7974 [00:02<00:00, 3296.79it/s]

```
1 texto_noticias["sent_polarity_str"].hist(bins=80)
```



```
1 texto_noticias
```

	texto_eng	label	set	char_len	normaliza	clean_text	token_len_limpio	sent_subjectivity_str	sent_polarity_str
0	[CHINA: releases a pandemic USA: releases UFO...	0	train	92	[china, release, pandemic, usa, release, ufos,...	china release pandemic usa release ufos argent...	10	0.000000	0.000000
1	San Francisco (USA) 100 years ago after experi...	0	train	95	[san, francisco, usa, 100, year, ago, experien...	san francisco usa 100 year ago experience 1918...	9	0.000000	0.000000
2	Please put on your mask or we will have 28 new...	0	train	93	[please, put, mask, 28, new, song, artist, col...	please put mask 28 new song artist collaborate...	9	0.454545	0.136364
3	The new „covid norms and reception rules“ in t... If confinement	0	train	210	[new, covid, norm, reception, rule, spanish, r...	new covid norm reception rule spanish retail t...	19	0.389773	0.105682

▼ 4. Vectorización

```
1 import gensim.downloader as api
2 glove_emb = api.load('glove-twitter-25') # Descargamos y cargamos el embedding de "glove-twitter-25"

1 def get_average_vector(sentence):
2     #retokenizamos con nuestra función
3     tokens = tokenize(sentence)
4     # Generamos lista de salida vacía
5     lista = list()
6     # Iteramos por cada token de la frase de entrada
7     for i in tokens:
8         # Si el token se encuentra en el embedding, añadir a la lista.
9         # Si no se encuentra (except), pasa al siguiente elemento.
10        try:
11            lista.append(glove_emb.get_vector(i))
12        except:
13            continue
14
15 # Calculamos el valor medio de los vectores generados
16 try:
17     resultado = np.mean(lista, axis=0) # 1vector - Dimension 25d
18 except:
19     # Si la lista está vacía, generar vector de ceros de tamaño el embedding
20     resultado = np.zeros(25)
21 return resultado

1 texto_noticias["embeddings"] = texto_noticias["clean_text"].progress_apply(lambda x: get_average_vector(x))
```

100%

7974/7974 [00:02<00:00, 3980.87it/s]

```

1 vector_data = pd.concat([texto_noticias.embeddings.apply(pd.Series),
2                           texto_noticias[["sent_polarity_str", "sent_subjectivity_str"]], axis=1)

1 vector_data.shape

(7974, 27)

1 vector_data = vector_data.fillna(0)

1 import scipy as sp
2 #Extraemos las etiquetas y las asignamos a la variable y
3 y = texto_noticias["label"].values.astype(np.float32)
4 X = sp.sparse.csc_matrix(vector_data)

```

▼ 5. Entrenamiento y evaluación de modelos

```

1 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
2 print(X_train.shape)
3 print(X_test.shape)

(5980, 27)
(1994, 27)

```

Rebalanceo

```

1 from collections import Counter
2
3 from sklearn.datasets import make_classification
4
5 from imblearn.over_sampling import SMOTE
6
7 print('Original dataset shape %s' % Counter(y_train))
8
9 sm = SMOTE(random_state=42)
10
11 X_res, y_res = sm.fit_resample(X_train, y_train)
12
13 print('Resampled dataset shape %s' % Counter(y_res))
14

Original dataset shape Counter({0.0: 4586, 1.0: 1394})
Resampled dataset shape Counter({1.0: 4586, 0.0: 4586})

1 np.array(np.unique(y_res, return_counts=True)).T

array([[0.000e+00, 4.586e+03],
       [1.000e+00, 4.586e+03]])

```

Entrenamiento

```

1 # Definimos las funcionalidades pertinentes de sklearn:
2 from sklearn.model_selection import KFold
3 from sklearn.model_selection import cross_val_score
4 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
5 import warnings
6 # Definimos la función encargada de evaluar los modelos:
7 def model_evaluation(models, score, X, y):
8     results = []
9     names = []
10    #Para cada modelo
11    for name, model in models:
12        warnings.filterwarnings('ignore')
13        # Generamos un Kfold
14        KF = KFold(n_splits = 10, shuffle = True, random_state = 98)
15
16        # hacemos cross_val
17        cv_results = cross_val_score(model, X, y, cv = KF, scoring = score, verbose = False)
18
19        # Guardamos los resultados:
20        results.append(cv_results)
21        names.append(name)
22
23    # Mostramos los resultados numéricamente:

```

```

24     print('Metric: {} , KFold '.format(str(score)))
25     print("%s: %f (%f) " % (name, cv_results.mean(), cv_results.std()))
26
27     return results, names

```

```

1 # Cargamos los modelos
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.svm import SVC
6 from sklearn.model_selection import KFold
7 from sklearn.model_selection import GridSearchCV

1 # Definimos los modelos y generamos una lista con cada uno de ellos:
2 models = [
3     ("Logistic",LogisticRegression(random_state=30)),
4     ("SVC", SVC()),
5     ("DecisionTreeClassifier", DecisionTreeClassifier())
6 ]
7
8 evaluation_score = "accuracy"
9
10 model_evaluation(models, evaluation_score, X_res, y_res)

```

```

Metric: accuracy , KFold
Logistic: 0.535101 (0.019627)
Metric: accuracy , KFold
SVC: 0.594092 (0.012386)
Metric: accuracy , KFold
DecisionTreeClassifier: 0.694290 (0.014047)
([array([0.56535948, 0.55882353, 0.52126499, 0.56161396, 0.53544166,
         0.50817884, 0.51145038, 0.53762268, 0.52344602, 0.52780807]),
 array([0.5751634 , 0.59803922, 0.58015267, 0.58778626, 0.608506 ,
         0.59869138, 0.61286805, 0.58015267, 0.60632497, 0.59323882]),
 array([0.67973856, 0.68300654, 0.70338059, 0.70883315, 0.70010905,
         0.7044711 , 0.69138495, 0.66303162, 0.70665213, 0.70229008])),
 ['Logistic', 'SVC', 'DecisionTreeClassifier'])

```

Vamos a realizar un GridSearch para encontrar los mejores parametros del Decision Tree Classifier.

```

1 model = DecisionTreeClassifier()
2 grid_param = {
3     'criterion': ['gini', 'entropy'],
4     'max_depth': range(1,10),
5     'min_samples_split': range(1,10),
6     "min_samples_leaf": range(1,5)
7 }

1 model_grid = GridSearchCV(estimator=model,
2                             param_grid=grid_param,
3                             scoring='accuracy',
4                             cv=10,
5                             n_jobs=-1,
6                             error_score=0)

1 grid_result = model_grid.fit(X_res, y_res)

1 print(grid_result.best_params_)

{'criterion': 'gini', 'max_depth': 9, 'min_samples_leaf': 2, 'min_samples_split': 4}

1 print(grid_result.best_score_)

0.5975905374872594

1 from sklearn.model_selection import (KFold, cross_val_score,cross_validate)
2 from sklearn.metrics import make_scorer, accuracy_score, precision_score, recall_score, f1_score
3
4 model=DecisionTreeClassifier(criterion= 'gini',
5                             max_depth= 9,
6                             min_samples_split= 4,
7                             min_samples_leaf= 2)
8 model.fit(X_res,y_res)

DecisionTreeClassifier(max_depth=9, min_samples_leaf=2, min_samples_split=4)

```

```

1 from sklearn.metrics import f1_score, confusion_matrix
2 from sklearn.metrics import classification_report
3 y_pred = model.predict(X_test)
4 print(confusion_matrix(y_test,y_pred))
5 print(classification_report(y_test,y_pred))

```

```

[[757 766]
 [253 218]]

```

	precision	recall	f1-score	support
0.0	0.75	0.50	0.60	1523
1.0	0.22	0.46	0.30	471
accuracy			0.49	1994
macro avg	0.49	0.48	0.45	1994
weighted avg	0.62	0.49	0.53	1994

Podemos observar como claramente clasifica mejor la clase 0 que la 1. No es un gran modelo, tiene un 49% de probabilidad de predecir correctamente. Realmente el modelo no funciona bien, en mi opinión he pre-procesado correctamente los datos con los que he vectorizado y he rebalanceado los datos, no entiendo los resultados. Después de muchas pruebas mis resultados son estos, tampoco sé si me tendría que salir un valor muy elevado, pero seguro que podría ser mejor.

- Es posible que al traducir el corpus haya tenido algún problema, ya que me resultó complicado.
- He obviado extraer el sentimiento de los emojis, quizá hubiese aportado un poco más de valor.
- Se que hay el Random Forest Classifier y el Naive Bayes Classifier que pueden tener buenos resultados en el NLP, igual podría aumentar el accuracy del F1-score.