	IEFI Lautaro Santos Da Silveira  Librerias
In [326	<pre>import numpy as np import nltk import nltk nltk.download('punkt') from nltk import regexp_tokenize from nltk.corpus import stopwords</pre>
	<pre>nltk.download('stopwords') !python -m spacy download es_core_news_sm import spacy nlp= spacy.load('es_core_news_sm') import re from sklearn.feature_extraction.text import TfidfVectorizer from sklearn.utils import shuffle</pre>
	from sklearn.model_selection import train_test_split from wordcloud import WordCloud import matplotlib.pyplot as plt from sklearn.model_selection import GridSearchCV from sklearn.svm import SVC
	<pre>[nltk_data] Downloading package punkt to [nltk_data] C:\Users\Usuario\AppData\Roaming\nltk_data [nltk_data] Package punkt is already up-to-date! [nltk_data] Downloading package stopwords to [nltk_data] C:\Users\Usuario\AppData\Roaming\nltk_data [nltk_data] Package stopwords is already up-to-date!</pre>
In [ ]:	Set de datos  #df= pd.read_excel(r'D:\Escritorio\IES21\4to cuatrimestre\PLN\Noticias Verdaderas y Falsas.xlsx')  df= pd.read_excel(r'C:\Users\Usuario\Desktop\lautaro\ies 21\IES21\4to cuatrimestre\PLN\Noticias Verdaderas y Falsas.xlsx')
Out[ ]:	
	<ul> <li>True Los mejores momentos del desfile de Victoria's</li> <li>True Los jóvenes de entre 19 y 24 años dedican meno</li> <li>True Los españoles que aguardaron horas delante de</li> <li>False Los eurodiputados de la comisión de Sanidad de</li> </ul>
	1816 False En busca y captura el joven que lanzó pirotecn  Divido mi Set de Datos en dos, generando uno que serán los casos que tengo conocimiento historico, perteneciente a df_conocido, y otro que no tendrán la clasificacion, los cuales se utilizarán en un futuro para la predicción.
In [ ]:	<pre># Set de datos para entrenamiento y elección de modelo df_conocido=df_conocido.reset_index() df_conocido.drop(['index'],axis=1,inplace=True) # Set de datos a predecir df_Testeo = df_desconocido['Text']</pre>
	<pre>print('Set de Datos conocidos') display(df_conocido.head()) print('Set de Datos desconocidos') df_Testeo.head()  Set de Datos conocidos class</pre>
	<ul> <li>False Matan a la madre, la descuartizan y aún con la</li> <li>True Tras cerrar con varios hilos abiertos JuegodeT</li> <li>True La candidata socialista a la presidencia de La</li> <li>False Pese a ser el animal doméstico por excelencia,</li> </ul>
Out[ ]:	4 True Aramco, la empresa más rentable del planeta,  Set de Datos desconocidos  1732 Un sacerdote español fue expulsado de su parro  370 BBVA ha reconocido a la Comisión Nacional del
	656 El Govern de la Generalitat responsabilizó aye 349 Arrestado un hombre de 24 años sospechoso de s 1920 España registró el mayor descenso interanual Name: Text, dtype: object  Gráfico de proporcion en mi Set de Datos conocidos
In [ ]:	Se desea analizar si tengo variacion en mi set de datos luego de hacer una separacion aleatoria.  y = np.array([(df_conocido['class']==True).sum(), (df_conocido['class']==False).sum()])  mylabels = ["Verdaderas", "Falsas"]  plt.pie(y, labels = mylabels)
	plt.show()  Verdaderas
	Falsas
In [ ]:	Funcion de preprocesamiento  stops= set(stopwords.words('spanish'))
	<pre>def preprocesamiento(documento):     # Eliminación de valores numéricos     documento = re.sub('\d', ' ', documento)     # Minusculizacion del documento     documento = documento.lower() ####################################</pre>
	# Tokenizo con una List Comprehension, y luego a esta le filtro las stopwords tokenizado = [token for token in regexp_tokenize(documento, pattern='\w+')] filtro = [filtro for filtro in tokenizado if filtro not in stops] ####################################
	# Lematizo al texto, para eliminar los signos de puntuación o demas signos que no me aportaran al analisis, se filtraran los tokens # que están en la anterior List Comprehension llamada filtro lemma = [lema.lemma_ for lema in doc if str(lema) in filtro] # transformo mi lista de tokens a una cadena de texto Procesada=" ".join(map(str, lemma)) return Procesada
In [ ]:	<pre>x= df_conocido['Text'] y=df_conocido['class']</pre>
Out[ ]:	pre_procesamiento  matar madre descuartizar aún carne caliente co  tras cerrar varios hilo abierto juegodetrono h  candidata socialista presidencia rioja conchab  pese ser animal doméstico excelencia españa ge
	aramco empresa rentable planeta lanzar domingo  1595 arabio saudito acabar cortar él cabeza esra ll 1596 maternidad frustrado españolasparda montar dec 1597 agente policía guardia civil trabajar programa 1598 acción buen buen forzar reacción violento inde 1599 ex cuñada josé enrique abuín gey alias chicle
In [ ]:	Name: Text, Length: 1600, dtype: object
Out[]:	df_vector = pd.DataFrame(vectores, columns=nombres ) df_vector  abalanzar abandonado abandonar abandono abascal abascal abateír abatir abc abcel órbita órgano óscar últimamente último único útil coro tapa te  0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
	1       0.0
	1598 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.
	Graficos Palabras mas usadas en todo el set de datos
In [ ]:	<pre>todas_palabras = " ".join([palabra for palabra in pre_procesamiento]) wc = WordCloud(background_color = "black",colormap = "hsv", max_font_size = 150, random_state = 123).generate(todas_palabras) plt.imshow(wc, interpolation = "bilinear") plt.axis("off") plt.show()</pre>
	partido barcelona de la contro cario con contro cario cario cario cario municipio cario cario municipio cario
	Indicate the second sec
	Palabras mas usadas en noticias verdaderas
In [ ]: In [ ]:	<pre>df_Grafico['Token']= pre_procesamiento df_Grafico['clase']=y</pre>
	plt.imshow(wc, interpolation = "bilinear") plt.axis("off") plt.show()
	menos in the pedir
	proximo proximo partido partido partido parte pa
In [ ]:	Palabras mas usadas en noticias falsas  Palabras_Falsas = " ".join([palabra for palabra in df_Grafico['Token'][df_Grafico['clase']==False]]) wc = WordCloud(background_color = "black", colormap = "hsv", max_font_size = 150, random_state = 123).generate(Palabras_Falsas)
	plt.imshow(wc, interpolation = "bilinear") plt.axis("off") plt.show()
	persona persona per per per persona per
	Cataluna de jar informar pedir infor
To [ ].	Modelos de ML  # Conoro mi cot do datos para entrepemiento y tostos
In [ ]: In [ ]: Out[ ]:	<pre>x_train,x_test,y_train,y_test= train_test_split(df_vector,y,test_size=0.2,random_state=123) len(x_train), len(x_test)</pre>
In [ ]:	Modelo de Clasificacion SVM
	param_grid_svc = [{'C':[1.32,1.33,1.34], 'gamma':['scale']}]  # Creamos el GSCV para buscar los mejores hiperparámetros gscv_svc = GridSearchCV(estimator=svc,
	n_jobs=-1, cv=5, verbose=0, refit=True)  gscv_svc.fit(x_train, y_train)
Out[ ]:	[
In [304 Out[304	gscv_svc.best_params_ {'C': 1.32, 'gamma': 'scale'}
In [330	<pre>desvio_mejor_svc=desvios_svc[gscv_svc.best_index_] print("Accuracy: ",gscv_svc.best_score_,' +/- ', 2*desvio_mejor_svc, '( 95% )') Accuracy: 0.771875 +/- 0.06434768838116875 ( 95% )</pre>
In [313	<pre>svc_final= SVC(random_state=123, C=1.33, gamma='scale') svc_final.fit(x_train, y_train) y_pred = svc_final.predict(x_test)</pre>
	AC = accuracy_score(y_test, y_pred)  # Matriz de confusión  conf_matrix = confusion_matrix(y_true=y_test, y_pred=y_pred)  fig, ax = plt.subplots(figsize=(5, 5))  ax.matshow(conf_matrix, cmap=plt.cm.Blues, alpha=0.3)  for i in range(conf_matrix_shape[0]):
	<pre>for i in range(conf_matrix.shape[0]):     for j in range(conf_matrix.shape[1]):         ax.text(x=j, y=i,s=conf_matrix[i, j], va='center', ha='center', size='xx-large')  plt.xlabel('Predictions', fontsize=18) plt.ylabel('Actuals', fontsize=18) plt.title('Confusion Matrix', fontsize=18)</pre>
	print(f'El valor del Accuracy del modelo es igual a {AC}') plt.show()  El valor del Accuracy del modelo es igual a 0.76875  Confusion Matrix  On Matrix  Confusion Matrix
	ο- 117 49 <u>«</u>
	Actuals
	1- 25 129
	Predictions  Mi set de testeo cuenta con 320 observaciones, y mi modelo predice correctamente 246 observaciones.
In [320	<pre>nuevos_transf= vectorizer.transform(nuevos).toarray()</pre>
In [321	<pre>svc_final.fit(df_vector,y) # Obtengo las predicciones de los datos nuevos predicciones=svc_final.predict(nuevos_transf)  c:\Users\Usuario\Desktop\lautaro\ies 21\entorno\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but SVC was fitted wi th feature names</pre>
	warnings.warn(