

Universidad Politécnica de Madrid

MÁSTER UNIVERSITARIO EN INTELIGENCIA ARTIFICIAL

# CLASIFICADOR DOCUMENTAL CON GLOSARIO

Asignatura: Ingeniería Lingüística

Autores:

Joel Pardo Ferrera

Daniel Carmona Pedrajas

Correos:

joel.pardof@alumnos.upm.es

daniel.carmonap@alumnos.upm.es

23 Diciembre 2022

# Índice

<b>1. Introducción</b>	<b>2</b>
1.1. Objetivo de la práctica . . . . .	2
1.2. Estructura . . . . .	2
<b>2. Obtención de textos</b>	<b>2</b>
2.1. Web scraping . . . . .	3
2.2. OCR . . . . .	3
<b>3. Preprocesamiento de datos</b>	<b>4</b>
3.1. Procesamiento de caracteres especiales . . . . .	4
3.2. Procesamiento de stop words . . . . .	4
3.3. Lematización . . . . .	4
3.4. División del conjunto de datos . . . . .	4
<b>4. Creación del glosario</b>	<b>5</b>
4.1. Método de creación del glosario . . . . .	5
4.2. Selección de términos para el glosario . . . . .	5
<b>5. Clasificación</b>	<b>6</b>
5.1. Clasificador mediante TF-IDF en Vector Space Model . . . . .	6
5.1.1. Elementos del Vector Space Model . . . . .	6
5.1.2. Proceso de clasificación . . . . .	6
5.1.3. Resultados . . . . .	7
5.2. Otros clasificadores . . . . .	7
5.2.1. Entradas . . . . .	8
5.2.2. Entrenamiento, evaluación y producción . . . . .	8
5.2.3. Regresión logística . . . . .	8
5.2.4. Naive Bayes . . . . .	9
5.2.5. SVM lineal . . . . .	9
5.2.6. SVM con kernel gaussiano . . . . .	9
5.2.7. Resultados . . . . .	9
<b>6. Conclusiones</b>	<b>11</b>
<b>A. Fuentes de noticias</b>	<b>11</b>
<b>B. Web scraping en detalle</b>	<b>13</b>
<b>C. OCR en detalle</b>	<b>14</b>
<b>D. Lista de stop words</b>	<b>17</b>
<b>E. Glosarios completos</b>	<b>18</b>
<b>F. Modelo tf-idf en detalle</b>	<b>21</b>
F.1. Proceso de clasificación . . . . .	21
F.2. Resultados en detalle de la clasificación . . . . .	22
<b>G. Otros clasificadores en detalle</b>	<b>26</b>

# 1. Introducción

En esta sección presentaremos los objetivos de la práctica y la estructura de la memoria.

## 1.1. Objetivo de la práctica

El objetivo principal de la práctica es construir un sistema de **clasificación de textos**. Este problema consiste en asignar una clase o categoría a un documento dada una descripción del mismo en un cierto espacio documental, es decir, una codificación donde esté comprimida la información que contiene [1].

La clasificación de textos se puede abordar con métodos de aprendizaje supervisado y no supervisado. En nuestro caso utilizaremos el **aprendizaje supervisado** ya que nuestros textos están etiquetados por su categoría. En total, contaremos con 240 textos o noticias que se reparten de la siguiente manera:

- 60 noticias sobre **deportes**. Para obtener una representación completa de la categoría deportes hemos recogido:
  - 15 noticias sobre fútbol.
  - 15 noticias sobre baloncesto.
  - 15 noticias sobre tenis.
  - 15 noticias sobre Fórmula 1.
- 60 noticias sobre **salud**.
- 60 noticias sobre **ciencia**.
- 60 noticias sobre **política**

Para crear nuestro clasificador usaremos 30 noticias de cada categoría y posteriormente lo validaremos con la mitad restante.

## 1.2. Estructura

En esta memoria se encuentran documentados los pasos que hemos seguido para construir nuestro clasificador de textos:

- En la sección 2 explicamos las técnicas utilizadas para **obtener** nuestros textos
- En la sección 4 explicamos la **creación del glosario** para cada categoría con el que crearemos nuestro modelo.
- En la sección 5 explicamos los **modelos** utilizados para la clasificación de textos y analizamos los resultados obtenidos.
- En la sección 6 comentamos las **conclusiones** finales de la práctica.
- Por último en el apéndice encontraremos tablas e información adicional.

# 2. Obtención de textos

Lo primero que necesitamos para clasificar textos son textos, podríamos haber optado por recolectar una serie de textos de alguna base de datos de fácil acceso que tenga los documentos que buscamos pero hemos decidido suponer que no tenemos acceso a estas herramientas para simular un problema real y aprender otros métodos que pueden ser útiles en un futuro: **Web scraping** y **OCR**.

Como hemos mencionado anteriormente, el número de noticias recogidas es de 240 en total, 60 por categoría, obteniéndolas de 65 fuentes distintas (ver fuentes en apéndice A) para evitar sesgos.

## 2.1. Web scraping

Web scraping [2] es una técnica que permite recoger información de varias páginas web de forma automática explotando lenguajes de marcado que se pueden encontrar en la web. Los pasos que hemos seguido han sido los siguientes:

1. Seleccionamos una categoría de noticias en *Google News*, véase deportes.
2. Obtenemos una lista de links a noticias de deportes.
3. Utilizando la librería de python *request* obtenemos el html de cada link.
4. En cada html usamos la librería de python *BeautifulSoup* y recogemos todo el texto contenido en el tag `<article>` de ese mismo html.
5. Guardamos un archivo *.txt* con el texto de la noticia en el directorio correspondiente.

Hemos usado esta técnica para obtener noticias sobre **deportes** y **salud**. Cabe destacar que las noticias obtenidas con este método son todas correspondientes al mismo día, creemos que esto puede afectar al sesgo de los datos y sería interesante recoger noticias que cubriesen un espacio de tiempo mayor, incluso de años para obtener un modelo que generalice mejor, pero este no es el objeto de la práctica.

El proceso de Web scraping se describe con más detalle en el apéndice [B](#).

## 2.2. OCR

Las siglas de esta técnica hacen referencia al reconocimiento óptico de caracteres, de un determinado alfaéto, dentro de una imagen. Originalmente, se utiliza para digitalizar documentos de forma efectiva, tanto temporal como computacionalmente.

Uno de los problemas principales que muestra es el mal reconocimiento de los caracteres que localiza dentro de la imagen. Este hecho se produce sobretodo cuando la calidad de la misma no es muy buena, cuando aparecen manchas sobre el texto, o cuando son textos manuscritos a mano.

Para la implementación de la misma hemos seguido los siguientes pasos:

1. Abrimos desde el navegador un portal de noticias. Buscamos las categorías en cuestión.
2. Seleccionamos una noticia y hacemos capturas de pantalla de la misma.
3. Almacenamos de forma ordenada las capturas de cada noticia correspondiente a una categoría.
4. Para cada noticia dentro de cada categoría, cargamos las capturas realizadas con la librería de python *PIL*.
5. Aplicamos el OCR a partir de la librería de python *pytesseract* especificando el lenguaje castellano.
6. Guardamos el resultado del OCR en un archivo en formato *.txt*.

Esta técnica únicamente la hemos empleado para obtener las noticias de las categorías **ciencia** y **política**. Como bien ocurre con la técnica explicada anteriormente, las noticias están fechadas en el mismo día, de recopilación de las mismas. Además de afectar al sesgo de los datos, el clasificador corre el riesgo de aprender el vocabulario contemporáneo desestimando así palabras que en un pasado no muy lejando podrían haber sido empleadas para la redacción de esa misma noticia.

La obtención de noticias mediante OCR se describe con más detalle en el apéndice [C](#).

### 3. Preprocesamiento de datos

En esta sección presentamos cómo hemos normalizado nuestros documentos. En este punto, hemos obtenido 240 textos mediante Web Scraping y OCR pero su contenido está en crudo, es decir, contiene signos de puntuación, tildes, saltos de línea y otros caracteres que no nos interesan a la hora de clasificar nuestro texto mediante un glosario. Aunque estos elementos no nos son útiles ahora, sí que podrían ser importantes si estuviésemos extrayendo conocimiento del texto, por ejemplo, la presencia de una interrogación puede preceder a una sección de texto con mucha información.

Además de esto eliminaremos las stop words y obtendremos el lema de la palabra.

Se puede ver toda una traza de ejecución del preprocesado en el notebook [Preprocessing.ipynb](#).

#### 3.1. Procesamiento de caracteres especiales

Como hemos comentado, vamos a prescindir de caracteres especiales así que primero quitamos los caracteres no ASCII, acentos y letras mayúsculas recurriendo a la librería de python *unicode* y eliminamos los signos de puntuación mediante *spaCy*. Un ejemplo de conversión sería el siguiente:

'El fútbol es un deporte.' → 'el futbol es un deporte'

La razón por la que decidimos eliminar los acentos es porque son unos de los principales errores ortográficos que se cometen y podríamos obtener frecuencias de palabras incorrectas perdiendo así información.

#### 3.2. Procesamiento de stop words

Gracias a la librería de python *spaCy* podemos usar un modelo del lenguaje español muy completo con el que podemos parsear, lematizar e incluso obtener análisis morfológicos de palabras. También cuenta con una lista de stop words de 521 palabras que podemos encontrar en el apéndice [D](#).

Las stop words son palabras muy frecuentes de la lengua que no suelen aportar información útil sobre el documento al que pertenecen, es por esto que decidimos prescindir de ellas. De esta manera nuestra frase anterior se convertía en:

'el futbol es un deporte' → 'futbol deporte'

#### 3.3. Lematización

El último paso que damos para terminar de procesar nuestros textos es obtener el lema de los términos que no han sido filtrados hasta este punto. Lematizar un término consiste en encontrar el mayor representante de una palabra y sus derivaciones morfológicas. La librería *spaCy* cuenta con un modelo del lenguaje español llamado *es\_core\_news\_lg* con el que podemos realizar esta tarea. Así tendríamos una conversión final como esta:

'chicago sufriendo' → 'chicago sufrir'

Usando los lemas de los términos agrupamos todas las derivaciones morfológicas de una palabra en un mismo elemento, dando así más peso a este grupo de palabras que nos puede dar información sobre la categoría a la que pertenece un texto. Por ejemplo, si en un texto aparecen los términos 'marcó', 'marcaba' y 'marcaron' únicamente contaríamos una repetición por cada uno de ellos, por lo que la información se dispersa, si en cambio usamos su lema 'marcar' tenemos una frecuencia de tres, lo que nos ayuda a detectar que el texto al que pertenecen estos términos posiblemente pertenezca a la categoría de deportes.

#### 3.4. División del conjunto de datos

Como sabemos estamos en el proceso de crear un clasificador de textos mediante aprendizaje supervisado, así que para crear nuestra herramienta vamos a utilizar un conjunto de entrenamiento y para validarlo un conjunto de validación por lo que en este momento vamos a dividir nuestras noticias en dos conjuntos:

- 120 noticias para crear nuestro clasificador guardadas en `noticias_train.csv`
- 120 noticias para validar nuestro clasificador guardadas en `noticias_test.csv`

En cada conjunto hemos introducido de forma aleatoria 30 noticias de cada categoría para mantener nuestros conjuntos balanceados.

## 4. Creación del glosario

En esta sección explicaremos cómo hemos conseguido nuestro glosario. Para detalles de la implementación consultar el notebook [Glosarios.ipynb](#) donde se muestra la traza de la última ejecución.

En primer lugar, debemos definir qué es un glosario. Pues bien, un *glosario* es un conjunto de términos representativos de un conjunto de documentos que permite agruparlos bajo una misma categoría. En nuestro caso, queremos obtener 4 glosarios distintos:

- Glosario de **deportes**.
- Glosario de **salud**.
- Glosario de **ciencia**.
- Glosario de **política**

En la sección 5 utilizaremos estos glosarios para crear nuestro clasificador.

### 4.1. Método de creación del glosario

Para crear nuestro glosario usaremos el método **tf-idf**. El esquema **tf-idf** [1] es una combinación de las métricas *term frequency* e *inverse document frequency*:

- Mediante **term frequency** obtenemos la frecuencia en la que el término  $t$  aparece en el documento  $d$  y lo denotamos como  $tf_{t,d}$ .
- Mediante **inverse term frequency** obtenemos una medida de cómo de frecuente es un término  $t$  en un conjunto de documentos de tamaño  $N$  y lo denotamos como  $idf_t$ .

Cabe destacar que  $idf_t$  se puede expresar de muchas maneras, en nuestro caso usaremos la expresión  $\log \frac{N}{df_t}$  donde  $df_t$  es el número de documentos que contiene el término  $t$ .

Por lo tanto obtenemos la métrica **tf-idf**:

$$tf - idf_{t,d} = (1 + \log tf_{t,d}) \times idf_t = (1 + \log tf_{t,d}) \times \log \frac{N}{df_t} \quad (1)$$

En otras palabras,  $tf - idf_{t,d}$  asigna a cada término  $t$  un peso en un documento  $d$  tal que:

- El peso es alto cuando  $t$  aparece muchas veces en un número pequeño de documentos.
- El peso es medio cuando  $t$  aparece menos veces en un documento o aparece en muchos documentos.
- El peso es bajo cuando  $t$  aparece en muchos documentos.

### 4.2. Selección de términos para el glosario

Habiendo definido la métrica **tf-idf** de un término en un documento, seleccionamos los términos que pertenecen al glosario de un categoría de la siguiente manera:

1. Creamos un diccionario con todas las palabras de los documentos de la categoría en cuestión.
2. Calculamos el tf-idf de cada término en cada documentos y los sumamos obteniendo su peso tf-idf total.
3. Ordenamos los términos de mayor a menor según su peso tf-idf total.

4. Seleccionamos los 50 primeros términos.

Mostramos los primeros 5 términos de cada glosario en la tabla 4. Para comprobar los glosarios completos consultar el apéndice E.

deportes	salud	ciencia	política
seguidor	estigma	sal	silva
mbappe	atras	cola	rosell
correr	cancer	meteoro	vuelo
suarez	anemia	agujero	precio
resto	cafa	llama	torra
mans	aceite	gemanidas	reyes

Tabla 1: 5 primeros términos de cada glosario ordenados según peso tf-idf total

## 5. Clasificación

En esta sección mostraremos los distintos clasificadores y resultados obtenidos.

### 5.1. Clasificador mediante TF-IDF en Vector Space Model

Aprovechando que estamos representando nuestros documentos como vectores, vamos a utilizar el Vector Space Model que nos proporciona tf-idf para clasificar las noticias

**Vector Space Model** [1] es un espacio vectorial común a una serie de vectores que codifican documentos y que permite crear clusters de documentos, clasificar documentos y crear métricas para puntuar documentos en base a una query.

#### 5.1.1. Elementos del Vector Space Model

Los elementos que necesitamos para representar nuestros documentos en forma vectorial son:

- **Diccionario**: un conjunto de términos que nos permitan representar todos los documentos en base a ellos. En nuestro caso el diccionario es la **unión de términos de todos los glosarios** obtenidos en la sección 4
- Un método de **vectorizar** nuestros documentos en base a nuestro diccionario. En nuestro caso usaremos el método **bag of words** ,proporcionado por la librería *gensim*, que para cada documento guarda un vector con la frecuencia de cada palabra del diccionario.
- **Modelo tf-idf** creado con los **bag of words** de los documentos para ahora poder expresarlo con su *tf-idf* que como hemos visto anteriormente contiene mejor la información del documento.

#### 5.1.2. Proceso de clasificación

Para clasificar los documentos utilizaremos como referencia los glosarios de cada categoría.

Creamos un vector tf-idf de cada glosario y gracias a que estamos en un **Vector Space Model** [1] podemos compararlos con los vectores tf-idf de los documentos calculando su similaridad de similaridad del coseno:

$$sim(d1, d2) = \frac{d1 \cdot d2}{|d1||d2|} \quad (2)$$

Teniendo en cuenta esta métrica, calcularemos la similaridad entre el vector tf-idf de cada documento con los vectores de cada glosario y asumiremos que pertenece a la categoría con cuyo glosario tenga mayor similitud. Además de esto, guardamos los documentos en carpetas según la categoría asignada y ordenados de mayor a menor similitud con el glosario de la misma.

### 5.1.3. Resultados

De 120 documentos, hemos conseguido clasificar 111, los 9 no clasificados son documentos que han obtenido una similitud igual con todas las categorías o su representación en el modelo es un vector de 0s. Los resultados que obtenemos con respecto a los documentos clasificados son los siguientes:

Y la matriz de confusión:

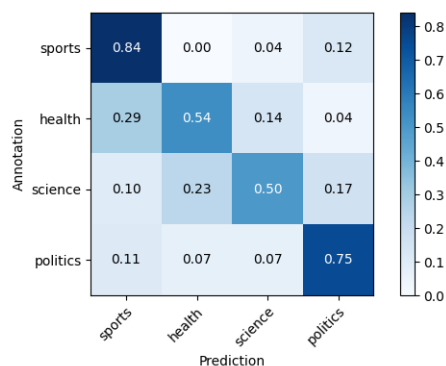


Figura 1: Matriz de confusión de nuestro clasificador

Hemos obtenido un modelo que no es demasiado bueno, como vemos en la matriz de confusión, las categorías con las que más problema tiene son con *health* y *science*, ambas tienen solapamiento semántico entre ellas porque la medicina pertenece al campo de la salud y la ciencia simultáneamente. Ocurre lo mismo con *health* y *sports*.

Ya habíamos comentado que nuestras noticias pertenecen a un periodo temporal muy pequeño y esto puede haber provocado que se centren simplemente en los temas de actualidad de ese momento, habiendo obtenido así un glosario sesgado y que no represente correctamente las características de cada una de las categorías, lo que podría explicar por qué ha habido documentos que no se han conseguido clasificar. Es por todo lo anterior que concluimos que para mejorar nuestros resultados con este modelo tendríamos que obtener un mejor glosario. De todas formas, decidimos probar con otro modelo para intentar mejorar nuestros resultados.

Para más detalle sobre la clasificación mediante tf-idf y resultados consultar el apéndice F y el notebook [Clasificacion-TF-IDF.ipynb](#) donde se encuentra la traza de la última ejecución realizada.

## 5.2. Otros clasificadores

Vamos a hablar de otros clasificadores como la regresión logística, el SVM lineal, SVM con kernel gaussiano y Naive Bayes. Todos ellos implementados en la librería *Scikit-Learn*.

Para la tokenización de palabras empleamos el modelo *es\_core\_news\_lg* de la librería *Spacy*. Este es un pipeline para invocar a un modelo optimizado para la CPU.

Hemos diseñado un proceso secuencial donde:

1. Probamos los clasificadores con las diferentes mutaciones entre los datos.
2. Escogemos los tipos de mutaciones de los datos que mejor rendimiento han obtenido junto con los cuatro clasificadores que han mostrado mejor rendimiento y variamos sus parámetros.
3. Escogemos los tres mejores modelos con los respectivos parámetros que mejor rendimiento aportan y los comparamos con la matriz de confusión.



### 5.2.1. Entradas

A partir vocabulario recogido para el modelo de la subsección anterior, generamos una bolsa de palabras. Invocamos también un modelo TF-IDF distinto.

Ahora, mediante la interacción de estos componentes junto con los corpus originales, vamos a generar mutaciones de las entradas. Finalmente, quedan cuatro tipos de entradas:

- Bolsa de palabras (BoW) a partir de los corpus originales.
- Características TF-IDF a partir de la bolsa de palabras.
- Promediado de los word embeddings de todos los corpus correspondientes a cada documento a partir de dichos corpus originales.
- Características TF-IDF con los word embeddings promediados a partir de la bolsa de palabras, de las características TF-IDF y de los corpus originales.

El que mejor funciona es el tercero; promediado de los word embeddings de todos los corpus correspondientes a cada documento a partir de dichos corpus originales.

### 5.2.2. Entrenamiento, evaluación y producción

El proceso de aprendizaje del modelo se lleva de forma secuencial; primero el entrenamiento para enseñarle al modelo las muestras de cada una de las clases, evaluación del entrenamiento para ver cómo de bien el modelo es capaz de clasificar y por último la producción donde empleamos el modelo para predecir muestras nuevas que nunca ha visto.

Como método de evaluación nos basamos tanto en métricas como en la matriz de confusión. Para recordar voy a hacer una síntesis de estos dos elementos.

La matriz de confusión es una matriz que permite la visualización directa del rendimiento de un algoritmo. Cada fila como cada columna representa una de las clases. La diferencia entre las filas y columnas es que las primeras son la clase verdadera mientras la segunda es la clase predicha.

Por ejemplo, tenemos un documento de la categoría deportes y lo clasificamos. El modelo puede predecir que ese documento es de deportes o es de cualquier otra categoría. El modelo clasifica dicho documento como deportes, por lo tanto ha acertado. El modelo clasifica dicha noticia como ciencia, por lo tanto ha fallado. Estamos, en definitiva, viendo el error del clasificador.

Las métricas numéricas diseñadas para evaluar el rendimiento del modelo. Están muy cercanas a la matriz de confusión. Las métricas que empleamos son las comunes:

- Exactitud (Accuracy)
- Precision
- Recall (Exhaustividad)
- F1 score

### 5.2.3. Regresión logística

Es un tipo de análisis donde la variable objetivo, es decir, la variable a predecir es una variable categórica a partir de unas variables predictoras (características). Una variable categórica es una variable que puede adoptar un número limitado de estados (categorías).

Este modelo es de los mejores en todas las iteraciones que hemos hecho. Ha llegado, en muchas ocasiones a la última fase y en varias, ha sido seleccionado al mejor modelo. En la segunda fase variamos el *solver*, es decir, el algoritmo que utiliza para optimizar el resultado para el problema dado. El mejor valor de este parámetro es *liblinear* y *newton-cholesky*.

#### 5.2.4. Naive Bayes

Este clasificador entra dentro de la familia de clasificadores probabilísticos. Se basa en el teorema de Bayes con cierta independencia (naive) entre las características.

El rendimiento de este modelo es bastante bueno. En la segunda aproximación cuando variamos el parámetro que modifica la varianza para calcular la estabilidad *var\_smoothing*, no provoca cambios en el rendimiento. En ninguna de las iteraciones que hemos hecho, ha llegado a la última fase seleccionado como uno de los mejores modelos.

#### 5.2.5. SVM lineal

Entra dentro la familia de clasificadores lineales como la regresión logística o las máquinas de vectores de soporte con un descenso del gradiente estocástico como sería el algoritmo de optimización SGD. En este caso, el gradiente es estimado para cada muestra en una época y el modelo se actualiza al final reduciendo la tasa de aprendizaje.

Sin duda, es de los que mejores resultados nos ha proporcionado. Ha llegado en todas las iteraciones a la última fase y ha sido en todas seleccionado de los mejores modelos. En la segunda fase modificamos la función de pérdida que usa *loss*. Los mejores valores de este parámetro son *squared\_hinge*, *modified\_huber*, *perceptron* y *hinge*, es decir, casi todos los valores disponibles.

#### 5.2.6. SVM con kernel gaussiano

Clasificador C-vector de soporte, es decir, una máquina de vectores de soporte gaussiana. En este caso, separa espacios de puntos para todas las categorías y luego comienza a esbozar fronteras de separación entre los distintos conjuntos intentando maximizar la distancia entre ellos. El vector de soporte se realiza con el esquema uno a uno, es decir, clase contra clase.

Los resultados de este modelo no son buenos. No pasa de la segunda fase donde variamos el parámetro *gamma* que varia el coeficiente del kernel.

#### 5.2.7. Resultados

El mejor modelo es son las SVM con kernel lineal mientras que el peor son las SVM con kernel gaussiano (no lineal).

Naive Bayes es bueno pero no lo suficiente. La regresión logística es el segundo mejor modelo. Introducimos los matrices de confusión de los tres primeros.

	Accuracy	Recall	Precision	F1 score
<b>modified_huber - Linear SVM</b>	88.47	86.7	88.47	88.4
<b>liblinear - Logistic Regression</b>	87.5	88.2	87.5	87.7
<b>newton-cholesky - Logistic Regression</b>	85.8	86.9	85.8	86.2
<b>log - Linear SVM</b>	85	87.7	85	85.1
<b>hinge - Linear SVM</b>	84.2	85.6	84.2	84.5
<b>perceptron - Linear SVM</b>	83.3	85.2	83.3	83.7

Tabla 2: Métricas de los mejores clasificadores en la tercera fase

Para más detalle sobre la clasificación mediante tf-idf y resultados consultar el apéndice [G](#) y el notebook [Clasificacion-Resto.ipynb](#) donde se encuentra las trazas de la última ejecución realizada.

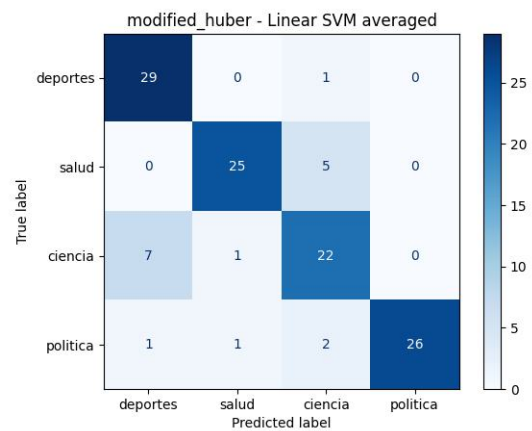


Figura 2: Matriz de confusión de modified\_huber - Linear SVM

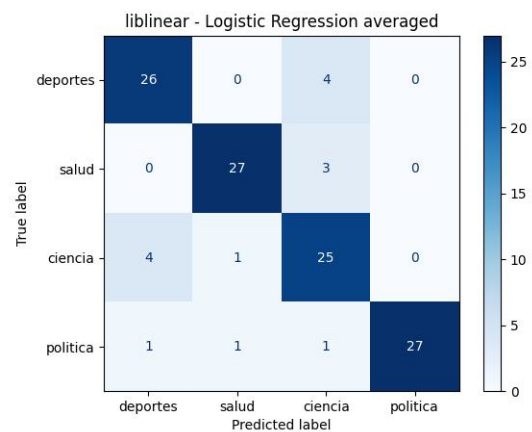


Figura 3: Matriz de confusión de liblinear - Logistic Regression

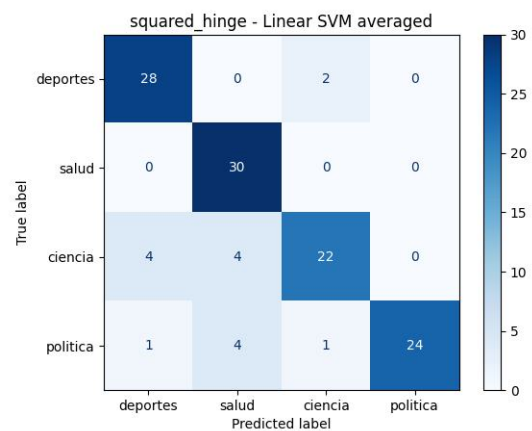


Figura 4: Matriz de confusión de squared\_hinge - Linear SVM

## 6. Conclusiones

Hemos tocado desde la obtención hasta la obtención de resultados de los clasificadores. Finalmente concluimos que:

- La obtención de noticias es más efectiva empleando la técnica de Web Scrapping. Sin embargo, para emplear el clasificador sobre noticias de fuentes analógicas, se podría utilizar de forma que agilizase el proceso de digitalización.
- En cuanto al preprocesado, hemos eliminado desde acentos hasta stop words. Hemos lematizado el texto y por último hemos dividido el conjunto de datos total en conjuntos definidos de entrenamiento y prueba.
- Para crear el glosario y poder acceder a él de forma sencilla, el método TF-IDF. De esta forma, seleccionamos las más relevantes por categoría que formarían parte de este conjunto de palabras.
- En cuanto a la clasificación, hemos planteado cinco modelos diferentes desde TF-IDF en Vector Space Model hasta una Regresión logística. El modelo que mejor rendimiento presenta de todos es un SVM con kernel lineal. Sin embargo, cabe destacar el buen rendimiento de la regresión logística posicionándose así como segundo en el ranking.
- Para finalizar, hemos cumplido con el objetivo de este proyecto que sería el de implementar un clasificador de textos. Hemos planteado varios y analizado sus rendimientos tras haber recopilado todos los datos necesarios y formado un glosario por cada una de las categorías que pretendíamos diferenciar.

### A. Fuentes de noticias

Las noticias recogidas vienen de 65 fuentes distintas:

- Antena 3 Noticias
- El Periódico
- El Confidencial
- ABC
- Okdiario
- Diario de Almería
- Geriatricarea.com
- Defensa Central
- laSexta
- Diario de Ávila
- Redacción Médica
- Planeta Triatlón
- elDiario.es
- MUI Kitchen
- Vista al Mar
- Periodismo.com
- ContraRéplica
- Newtral
- Car and Driver

- miarevista.es
- Fitness
- Cadena SER
- El Día de Valladolid
- Business Insider España
- ATP Tour
- Marca
- Medicina y Salud Publica
- ES
- AS
- 20minutos.es
- Faro de Vigo
- Superdeporte
- Eurosport ES
- Onda Cero
- Ser Padres
- El Debate
- El Imparcial de Oaxaca
- El Periódico de Extremadura
- Caracol Radio
- La Voz de Medina Digital
- Alimento
- NotiPress
- La Nueva España
- Murcia.com
- MSN
- asajaen.com
- Málaga Hoy
- La Razón
- La Voz de Galicia
- MARCA.com
- MUNDO DEPORTIVO
- Vitónica
- Diario de Castilla y León
- INFORMACIÓN

- Diario de Cádiz
- Infosalus
- Bernabéu Digital
- Qué!
- Noticias de Gipuzkoa
- EL ESPAÑOL
- TRT Español
- EL PAÍS
- nbamaniacs
- Deia
- Sporting News
- El Independiente
- El Mundo

## B. Web scraping en detalle

El código fuente que hemos utilizado para hacer Web scraping se encuentra en [google\\_news\\_url\\_scraping.py](#) y [text\\_scraping.py](#).

Lo primero que hacemos es obtener los links y títulos de las noticias mediante el script *google\_news\_url\_scraping.py*, ¿Cómo? Si accedemos a [Google News](#) y escogemos un tema o *topic* llegaremos a una página que tiene interfaz gráfica para el usuario, pero si modificamos la url y escribimos *rss* delante de *topic*, nos dirigiremos a una página como [esta](#):

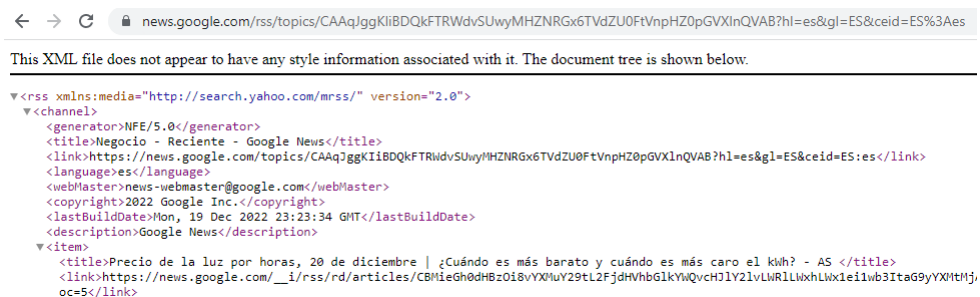


Figura 5: XML Google News

Es un archivo XML en el que cada <item> representa una noticia con su título y link que podemos parsear con la librería *BeautifulSoup* como se ve en el código [1](#).

```

1
2 # Obtenemos el archivo XML
3 response = requests.get(gnews_url)
4 # Parseamos el XML
5 soup = BeautifulSoup(response.text, "xml")
6 # Dividimos todos los items que contienen noticias
7 news = soup.find_all("item")
8 i = 0
9 retrieved = 0
10
11 # De cada item obtenemos el titulo y el link de la noticia

```

```

12 while retrieved < 15 and i < len(news)-1:
13     title = news[i].find('title')
14     link = news[i].find('link')
15     retrieved += 1
16     i += 1
17

```

Código 1: Uso de *BeautifulSoup* para scraping de links en Google News

Todos estos links se guardan en el dataframe *urls.csv* que tiene columnas de categoría, título y link con una fila por noticia.

Posteriormente con el script *text\_scraping.py* obtenemos el html haciendo un request para cada link y obteniendo el cuerpo de la noticia accediendo al cuerpo de la noticia que se encuentra dentro del tag `<article>` como se muestra en el código 2.

```

1 # Obtenemos el archivo html
2 response = requests.get(link)
3 # Parseamos el archivo html
4 soup = BeautifulSoup(response.text, 'html.parser')
5 # Nos quedamos con el tag article del html
6 article = soup.find('article')
7
8 if article == None:
9     print("Error, no se puede leer la noticia")
10    return None
11 # Eliminamos figuras
12 for fig in article.select('figure'):
13     fig.extract()
14 full_text = []
15 # Unimos todos los parrafos de article
16 for paragraph in article.find_all('p'):
17     full_text.append(paragraph.text)
18 '\n'.join(full_text)
19

```

Código 2: Guardado de noticias en documentos .txt

Hemos de aclarar que ha funcionado para la mayoría de links, aunque para algunos de ellos hemos tenido que eliminar a mano ciertas palabras que no pertenecían al artículo en sí como pueden ser hipervínculos o títulos de noticias relacionadas o incluso tener que copiar el texto del artículo directamente en el archivo *.txt*.

Por último, guardamos la noticia en su carpeta correspondiente como un archivo *.txt* como se ve en el fragmento de código 3.

```

1 def save_text(category, text, count):
2     file_name = f"{DATA_FOLDER}/{category}/{category}_{count}.txt"
3     if not exists(file_name):
4         with open(file_name, "w") as f:
5             if text is not None:
6                 f.write(text)
7

```

Código 3: Uso de *BeautifulSoup* para scraping de links en Google News

## C. OCR en detalle

El código fuente que hemos empleado para hacer OCR se encuentra en *OCR.py*, *registros.py* y *Carpetas.py*.

Comenzamos definiendo el espacio de almacenamiento. Cada carpeta categoría va a tener una carpeta. Cada subcarpeta de cada categoría es un número del 1 al 60, ya que 60 es el número máximo de noticias por categoría. Dentro de estas subcarpetas almacenamos las capturas de pantalla. Las imágenes con las que trabajamos están en formato *.png*.

Surge un problema a la hora de coger los títulos de las noticias. El OCR no distingue si el texto de la imagen que está visualizando se trata de un título o del cuerpo de una noticia. Por lo tanto, creamos otra subcarpeta

para almacenar las capturas de las imágenes que contienen los títulos. En esta carpeta tendremos 60 imágenes donde cada una representará el título de la noticia a la que hace referencia su nombre.

Por otra parte, para agrupar todos los documentos que extrae el OCR de una categoría concreta, hemos creado otra carpeta para este uso específico.

Empleamos la librería *os* para la creación de estos repositorios:

```
1 clases = ["Ciencia", "Politica"]
2
3 for z in clases:
4     # Introducimos el nombre de la clase y donde queremos que se guarden los documentos '.txt'
5     path = "../Documentos/" + z + "/"
6     path_guardado = path + "txt"
7     path_titulo = path + "titulo"
8
9     carpetas = os.listdir(path)
10
11     # Creamos los path de las carpetas (documentos)
12     for i in range(1, 61):
13         if str(i) not in carpetas:
14             os.mkdir(path + str(i))
15
16     if "txt" not in carpetas:
17         os.mkdir(path_guardado)
18
19     if "titulo" not in carpetas:
20         os.mkdir(path_titulo)
21
```

Código 4: Uso de *os* para la creación de los repositorios de almacenamiento

Vamos con la parte central de la cuestión a tratar en esta sección, el OCR. Su funcionamiento se basa en el recorrido de las subcarpetas de cada categoría creadas en el apartado anterior para obtener por cada una de ellas las capturas que forman una noticia y transformarla a un documento *.txt* almacenado en el directorio destinado a este fin. Podemos ver la forma en la que realiza esta tarea en el 5.

```
1 clases = ["Ciencia", "Politica"]
2
3 for z in clases:
4     # Introducimos el nombre de la clase y donde queremos que se guarden los documentos '.txt'
5     path = "../Documentos/" + z + "/"
6     path_guardado = path + "txt/"
7     path_titulo = path + "titulo/"
8
9
10    # Listamos todos los documentos para esta clase
11    carpetas = os.listdir(path)
12
13    # Eliminamos la carpeta de guardado y '.DS_Store'
14    carpetas = [c for c in carpetas if c != '.DS_Store' and c != 'txt' and c != 'titulo']
15
16
17    # Creamos las rutas de cada uno de los documentos
18    rutas = []
19    for i in carpetas:
20        rutas.append(path + i + "/")
21
22
23
24    # Para los documentos y las rutas
25    for c, j in zip(carpetas, rutas):
26
27        # Sacamos las imágenes por cada documento
28        imagenes = os.listdir(j)
29
30        # Copiamos el archivo que contiene un título a la carpeta de títulos
31        shutil.copy(j + "1.png", path_titulo + str(c) + ".png")
32
```



```

33     # Inicializamos lista a vacia. Sirve para guardar el texto de un mismo documento
34     texto = []
35
36     # Para cada una de las imagenes
37     for h in imagenes:
38         if h != '.DS_Store':
39             # Cargo la imagen
40             img = Image.open(j + h) # Abre la imagen con pillow
41             img.load()
42
43             # Saco el texto
44             text = pytesseract.image_to_string(img, lang='spa') # Extrae el texto de la
imagen
45             texto.append(text)
46
47
48     # Guardo el documento '.txt' en la carpeta de guardado
49     f = open(path_guardado+str(c)+".txt", 'w')
50     f.writelines(texto)
51     f.close()
52
53     print(c)
54
55
56
57     # Sacamos los documentos existentes en el path de guardado para comprobar que estan todos
58     documentos = os.listdir(path_guardado)
59     print(len(documentos))
60

```

Código 5: Uso de *pytesseract* y *PIL* para la creacion de los documentos *.txt* que representan cada una de las noticias

En una primera instancia hemos obtenido las noticias para dos de las categorías gracias al uso de Web Scrapping. Ahora que tenemos el conjunto de las noticias extraidas con OCR, vamos a generar un registro en formato *.csv* donde almacenamos información de cada una de las noticias. Nuevamente, se vuelve a emplear OCR sobre la carpeta de titulos para especificar este campo dentro del registro.

```

1  # Cargamos el registro
2  df = pd.read_excel("/Users/joelpardo/Desktop/Practica 2/Documentos/urls1.xlsx")
3
4
5  # Asignamos las variables
6  indices = df["Unnamed: 0"].tolist()
7  category = df["category"].tolist()
8  title = df["title"].tolist()
9  links = df["link"].tolist()
10
11
12  generico = "../Documentos"
13  path = generico + "/Links"
14
15
16  clases = ["Ciencia", "Politica"]
17  contador = indices[len(indices)-1] # 119
18
19  for z in clases:
20
21      doc = path + "/" + str(z) + '.txt'
22
23      path_titulos = generico + "/" + str(z) + "/titulo"
24
25
26      titulos = os.listdir(path_titulos)
27      titulos = [f for f in titulos if f != ".DS_Store"]
28
29
30      with open(doc, 'r') as f:
31          lin = f.readlines()
32          c = 0
33

```

```

34         for i in lin:
35             if i != "\n":
36                 if z == "Ciencia":
37                     clase = "science"
38                     category.append(clase)
39
40                 if z == "Politica":
41                     clase = "politics"
42                     category.append(clase)
43
44             print(c)
45             # Cargo la imagen
46             img = Image.open(path_titulos + "/" + titulos[c]) # Abre la imagen con pillow
47             img.load()
48
49
50             # Saco el texto
51             text = pytesseract.image_to_string(img, lang='spa') # Extrae el texto de la
imagen
52             text = text.replace('\n', '')
53             title.append(text.strip())
54             c += 1
55
56             new_string = i.replace('\n', '')
57             links.append(new_string.strip())
58
59
60             contador += 1
61             indices.append(contador)
62
63
64     print(len(indices))
65     print(len(category))
66     print(len(title))
67     print(len(links))
68
69
70     df_new = pd.DataFrame()
71     df_new["index"] = indices
72     df_new["category"] = category
73
74     path_rutas = "./Datos/Raw_data" #./TextClassification/
75
76     clases2 = df_new["category"].unique().tolist() #0-3
77
78     rutas_docs = []
79     r = []
80     for h in clases2:
81         for a in range(1,61):
82             p = path_rutas+ "/" + h + "/" + str(a) + ".txt"
83             rutas_docs.append(p)
84             r.append(a)
85
86     df_new["n_doc"] = r
87     df_new["title"] = title
88     df_new["path"] = rutas_docs
89     df_new["link"] = links
90
91     df_new.to_csv('../Documentos/urls.csv', header=True, index=False)
92

```

Código 6: Uso de *pandas* para la creacion del registro que representan cada una de las noticias

## D. Lista de stop words

La lista de stop words del modelo *es\_core\_news\_lg* de la librería *spaCy* es:

a, acuerdo, adelante, ademas, además, afirmó, agregó, ahí, ahora, ahí, al, algo, alguna, algunas, alguno, algunos, algún, allí, alí, alrededor, ambos, ante, anterior, antes, apenas, aproximadamente, aquel, aquella, aquellas, aquellos, aquí, aquél, aquélla, aquéllas, aquéllos, aquí, arriba, aseguró, así, así, atras, aun, aunque, añadió,

aún, bajo, bastante, bien, breve, buen, buena, buenas, bueno, buenos, cada, casi, cierta, ciertas, cierto, ciertos, cinco, claro, comentó, como, con, conmigo, conocer, conseguimos, conseguir, considera, consideró, consigo, consigue, consiguen, consigues, contigo, contra, creo, cual, cuales, cualquier, cuando, cuanta, cuantas, cuanto, cuantos, cuatro, cuenta, cuál, cuáles, cuándo, cuánta, cuántas, cuánto, cuántos, cómo, da, dado, dan, dar, de, debajo, debe, deben, debido, decir, dejó, del, delante, demasiado, demás, dentro, deprisa, desde, despacio, después, después, detras, detrás, día, días, dice, dicen, dicho, dieron, diez, diferente, diferentes, dijeron, dijo, dio, doce, donde, dos, durante, día, días, dónde, e, el, ella, ellas, ello, ellos, embargo, en, encima, encuentra, enfrente, enseguida, entonces, entre, era, eramos, eran, eras, eres, es, esa, esas, ese, eso, esos, esta, estaba, estaban, estado, estados, estais, estamos, estan, estar, estará, estas, este, esto, estos, estoy, estuvo, está, están, excepto, existe, existen, explicó, expresó, fin, final, fue, fuera, fueron, fui, fuimos, gran, grande, grandes, ha, haber, habia, habla, hablan, habrá, había, habían, hace, haceis, hacemos, hacen, hacer, hacerlo, haces, hacia, haciendo, hago, han, hasta, hay, haya, he, hecho, hemos, hicieron, hizo, hoy, hubo, igual, incluso, indicó, informo, informó, ir, junto, la, lado, largo, las, le, les, llegó, lleva, llevar, lo, los, luego, mal, manera, manifestó, mas, mayor, me, mediante, medio, mejor, mencionó, menos, menudo, mi, mía, mías, mientras, mio, míos, mis, misma, mismas, mismo, mismos, modo, mucha, muchas, mucho, muchos, muy, más, mí, mía, mías, mío, míos, nada, nadie, ni, ninguna, ningunas, ninguno, ningunos, ningún, no, nos, nosotras, nosotros, nuestra, nuestras, nuestro, nuestros, nueva, nuevas, nueve, nuevo, nuevos, nunca, o, ocho, once, os, otra, otras, otro, otros, para, parece, parte, partir, pasada, pasado, país, peor, pero, pesar, poca, pocas, poco, pocos, podeis, podemos, poder, podría, podriais, podriamos, podrian, podrias, podrá, podrán, podría, podrían, poner, por, porque, posible, primer, primera, primero, primeros, pronto, propia, propias, propio, propios, próximo, próximo, próximos, pudo, pueda, puede, pueden, puedo, pues, que, que, quedó, queremos, quien, quienes, quiere, quiza, quizás, quizá, quizás, quién, quiénes, qué, realizado, realizar, realizó, repente, respecto, sabe, sabeis, sabemos, saben, saber, sabes, salvo, se, sea, sean, segun, segunda, segundo, según, seis, ser, sera, será, serán, sería, señaló, si, sido, siempre, siendo, siete, sigue, siguiente, sin, sino, sobre, sois, sola, solamente, solas, solo, solos, somos, son, soy, su, supuesto, sus, suya, suyas, suyo, suyos, sé, sí, sólo, tal, tambien, también, tampoco, tan, tanto, tarde, te, temprano, tendrá, tendrán, teneis, tenemos, tener, tenga, tengo, tenido, tenía, tercera, tercero, ti, tiene, tienen, toda, todas, todavía, todavía, todo, todos, total, tras, trata, través, tres, tu, tus, tuvo, tuya, tuyas, tuyo, tuyos, tú, u, ultimo, un, una, unas, uno, unos, usa, usais, usamos, usan, usar, usas, uso, usted, ustedes, va, vais, vamos, van, varias, varios, vaya, veces, ver, verdad, verdadera, verdadero, vez, vosotras, vosotros, voy, vuestra, vuestras, vuestro, vuestros, y, ya, yo, él, ésa, ésas, ése, ésos, ésta, éstas, éste, éstos, última, últimas, último, últimos

## E. Glosarios completos

Los glosarios de train y test que contienen 50 términos por categoría son los siguientes:

<b>deportes</b>	<b>salud</b>	<b>ciencia</b>	<b>política</b>
seguidor	estigma	sal	silva
mbappe	atras	cola	rosell
correr	cancer	meteoro	vuelo
suarez	anemia	agujero	precio
resto	cafa	llama	torra
mans	aceite	gemanidas	reyes
boston	oliva	banyoles	plataforma
estabilidad	mascarilla	neandertal	anunciara
horford	pulmonar	congelacia	perao
cristiano	unicef	leo	sinema
exencion	cannabis	latigo	baron
enrique	pet	supersa	turismo
booker	bronquiolitis	cernan	ayres
tatum	requerir	quipus	animal
cambio	congelado	congelar	versia
formato	fisioterapia	nudo	podem
club	quiriaorgica	fuego	unilateral
juventus	creatividad	estampido	junts
ronaldo	tc	sonido	aborto
wiggins	rehabilitacia	menta	juventud
gonzalez	research	supermasivo	adjudicacia
gonzalo	signo	observacia	procesado
colombia	uk	perseidas	ex
magnussen	vejiga	inca	dema
steiner	marfan	matematica	magistrado
gasto	sindrome	alcoholismo	abascal
resistencia	deberiar	fumador	cgp
lakers	oxigeno	fumar	arizona
pts	lobo	izquierdo	desbloqueo
reb	manada	mutacion	perro
verdasco	parasito	coruaa	normalidad
kosmos	beneficio	diplococus	sandro
ktm	permafrost	roth	ribera
moto	recoletas	iphone	ven
atletico	mascaras	mandabula	tope
gimenez	fibrotico	martanez	geno
siebert	estancamiento	pertenecia	hidra
rez	expreso	sapiens	corrupto
dolares	viatris	galaxia	consistorio
grada	macula	calcio	bolaaos
krack	vision	punte	infraestructura
marko	farmaca	intercalar	ceuta
aerodinamico	utico	rotacia	deu
motogp	ir	atma	marruecos
golden	poveda	sfera	lamite
state	cereal	comercio	mejoraa
warriors	oleocantal	dinosaurio	aentiendena
medicamento	dilm	bunsen	denominador
diciembre	velazquez	comandante	entiendena
gira	gondii	masa	page

Tabla 3: Términos de cada glosario ordenados según peso tf-idf total

deportes	salud	ciencia	política
seguidor	arterial	vacuna	rosell
mbappe	strep	invierno	vuelo
warren	antidepresivo	gripe	insulto
doncic	bronquiolitis	navidad	junts
estabilidad	aprovechar	cuantico	precio
horford	temporada	captura	renovar
shapovalov	folcodina	honor	reyes
cristiano	depresion	fusia	rufian
exencion	cafa	cola	consulta
correr	mascarilla	gemanidas	anunciara
carpena	aceite	alcohol	sandro
enrique	oliva	latigo	monta
luis	entrenar	supersa	multilateral
mans	cabello	regalo	turismo
stakhovsky	estigma	tabaco	partida
malaga	epilepsia	quipus	preguntar
medico	quiriaorgica	agujero	ucraniano
booker	congelado	apple	ndum
aliassime	cuadro	antageno	refera
auger	gotlib	subtipo	nombramiento
pts	marfan	hemofilia	ribera
reb	manada	estampido	falso
nets	parasito	sonido	juventud
cantidad	colapsar	ue	podem
juventus	intensivo	km	unilateral
ronaldo	moralesla	fasica	urgente
clippers	padres	galaxia	educativo
lebron	torre	magic	auxiliar
garden	pez	reencuentro	lamite
florentino	comida	alineacia	normalidad
heredero	mascaras	almanaque	desbloqueo
gasto	permafrost	escarcha	ministros
atleta	ir	helada	designacia
colombia	poveda	tv	peticia
roger	amigdala	barrera	sortear
verdasco	hipocampo	mola	hipoteca
dificil	tpbm	perseidas	concejal
dolares	viatris	litio	ex
ktm	expreso	carbono	midcat
moto	crisis	groenlandia	diccionario
pa	saciant	sedimento	subjetivo
rez	vitonica	amerizaje	republicano
gimenez	flema	cuerda	actor
siebert	gondii	matematica	origen
uruguayo	gris	reactor	ceuta
collins	lader	cula	deu
rusia	intestinal	supermasivo	independencia
marko	anpe	vehaculo	tendrait
euros	masticar	fraa	elevado
marbella	laser	virus	mwh

Tabla 4: Términos de cada glosario ordenados según peso tf-idf total

## F. Modelo tf-idf en detalle

### F.1. Proceso de clasificación

Utilizamos la librería *gensim* de python para crear nuestros elementos del Vector Space Model.

En primer en el código ?? mostramos la creación de nuestro diccionario en base a nuestros glosarios.

```
1 def create_dictionary(glosarios):
2     # Obtenemos los tokens de nuestros glosarios
3     doc_tokens = [[termino for termino in glosario] for glosario in glosarios.values()]
4     # Los guardamos en un diccionario de la libreria gensim
5     dictionary = corpora.Dictionary(doc_tokens)
6     return dictionary
7
```

Código 7: Obtención del diccionario

En el código 8 mostramos la creación de los *bag of words* de nuestros documentos en base a nuestro diccionario:

```
1 def create_bag_of_words(docs_list, dictionary):
2     # Para cada documento creamos una lista de tokens
3     doc_tokens = [simple_preprocess(corpus) for corpus in docs_list]
4     # Para cada documento creamos un bow en base a nuestro diccionario
5     docs_bow = [dictionary.doc2bow(doc) for doc in doc_tokens]
6     return docs_bow
7
```

Código 8: Creación de bag of words

El código 9 muestra la conversión de nuestros documentos de *bag of words* a tf-idf

```
1 def create_tfidf(docs_bow):
2     docs_tfidf = models.TfidfModel(docs_bow, smartirs="lfc")
3     return docs_tfidf
4
```

Código 9: Conversión a tf-idf

Una vez realizadas las operaciones anteriores obtenemos las representaciones tf-idf de los glosarios como mostramos en el código 11

```
1 deportes_bow = glosarios_dict.doc2bow(glosarios["deportes"])
2 deportes_tfidf = tfidf_model[deportes_bow]
3
4 salud_bow = glosarios_dict.doc2bow(glosarios["salud"])
5 salud_tfidf = tfidf_model[salud_bow]
6
7 ciencia_bow = glosarios_dict.doc2bow(glosarios["ciencia"])
8 ciencia_tfidf = tfidf_model[ciencia_bow]
9
10 politica_bow = glosarios_dict.doc2bow(glosarios["politica"])
11 politica_tfidf = tfidf_model[politica_bow]
12
```

Código 10: Creación de bag of words

Calculamos la similaridad de cada documento con cada glosario y definimos su clase para guardarlos en un csv.

```
1
2 # Funcion que compara la similaridad de un documento con cada glosario y le asigna la
3   categoria con la que tenga mayor similitud
4 def get_prediction(n_documento,
5                   similarities_deportes,
6                   similarities_salud,
7                   similarities_ciencia,
8                   similarities_politica):
9     mejor_match = 0
10    mejor = "No clasificado"
11    if (similarities_deportes[n_documento] > mejor_match):
12        mejor_match = similarities_deportes[index]
```

```

12     mejor = "sports"
13     if (similarities_salud[n_documento] > mejor_match):
14         mejor_match = similarities_salud[index]
15         mejor = "health"
16     if (similarities_ciencia[n_documento] > mejor_match):
17         mejor_match = similarities_ciencia[index]
18         mejor = "science"
19     if (similarities_politica[n_documento] > mejor_match):
20         mejor_match = similarities_politica[index]
21         mejor = "politics"
22     return mejor, mejor_match
23
24 # Creamos las matrices de similaridad de gensim para cada glosario
25 index = similarities.MatrixSimilarity(tfidf_model[docs_bow])
26 similarities_deportes = index[deportes_tfidf]
27 similarities_salud = index[salud_tfidf]
28 similarities_ciencia = index[ciencia_tfidf]
29 similarities_politica = index[politica_tfidf]
30
31 # Para cada documento, definimos su categoria
32 predictions = []
33 similarities = []
34 for i in range(len(docs_bow)):
35     pred, sim = get_prediction(i,
36                               similarities_deportes,
37                               similarities_salud,
38                               similarities_ciencia,
39                               similarities_politica)
40     similarities.append(sim)
41     predictions.append(pred)
42
43 # Guardamos los resultados en un dataframe
44 noticias_test_dataframe.insert(2, "preds", predictions)
45 noticias_test_dataframe.insert(3, "similarity", similarities)
46

```

Código 11: Clasificación de los documentos

## F.2. Resultados en detalle de la clasificación

A continuación mostramos las noticias clasificadas en cada categoría ordenadas por similitud con la misma.

Título	Clasificación	Categoría original	Similaridad
Golden State asesta otro ...	sports	sports	0.356
Gobert toca fondo - AS ...	sports	sports	0.355
El fútbol es un cuento - ...	sports	sports	0.325
El fútbol fue lo de menos...	sports	sports	0.305
De jefe de estrategia de ...	sports	sports	0.273
Ferrari se enreda otra ve...	sports	sports	0.272
LeBron James: "No necesit...	sports	sports	0.272
Bolaños asegura que los b...	sports	politics	0.232
Domenicali: "Fernando Alo...	sports	sports	0.219
Djokovic, espectador de l...	sports	sports	0.219
Luis Enrique: «Si a Gavi ...	sports	sports	0.2
Los diagnósticos ocultos ...	sports	health	0.195
Houston da la sorpresa re...	sports	sports	0.194
El mejor golpe de Rafa Na...	sports	sports	0.187
¿Debo entrenar si estoy e...	sports	health	0.175
Sacramento arrasa en su v...	sports	sports	0.161
Rafa Nadal podría empezar...	sports	sports	0.155
Hawks y Nets han negociad...	sports	sports	0.155
Jimmy Butler regresará es...	sports	sports	0.155
El corredor de hidrógeno ...	sports	politics	0.153
UNAM se suma a la lucha c...	sports	health	0.152
T.J. Warren vuelve vuelve...	sports	sports	0.129
Los antidepresivos pueden...	sports	health	0.117
Europa recomienda retirar...	sports	health	0.117
Qué medicamentos tomar y ...	sports	health	0.105
. <sup>Es</sup> el mismo Alonso que c...	sports	sports	0.098
¿Qué trastorno se esconde...	sports	health	0.098
El Código Penal incluirá ...	sports	politics	0.078
La Agencia EspacialEspañ...	sports	science	0.078
Canadá Logra Su Primera C...	sports	sports	0.078
En todos los idiomas, las...	sports	science	0.077
Los pasajeros olvidados d...	sports	science	0.077
OMS cambia el nombre de l...	sports	health	0.072
Suspendido por una tangan...	sports	sports	0.072
Ferrero: "Trabajar con Zv...	sports	sports	0.058

Tabla 5: Documentos clasificados en la categoría *sports* ordenados por similaridad con el glosario



Título	Clasificación	Categoría original	Similaridad
Huevos de codorniz, una s...	health	health	0.297
“No se puede condenar a u...	health	health	0.258
El número de hogares acog...	health	politics	0.246
La revolucionariahipótesi...	health	science	0.225
¿De qué color es el fuego...	health	science	0.221
Los animales también tien...	health	science	0.218
Un fósil ubica al primere...	health	science	0.215
El aparato más desconocid...	health	science	0.209
¿Cómo prevenir la caída d...	health	health	0.191
¿Es cierto que comer rápi...	health	health	0.188
”Tenemos una epidemia de ...	health	health	0.174
BRONQUIOLITIS — Los pedia...	health	health	0.166
¿La bronquiolitis puede d...	health	health	0.166
Coosalud EPS conmemora el...	health	health	0.166
Conoce algunos de los rem...	health	health	0.166
La bronquiolitis. Pregunta...	health	health	0.145
Guía para entender la epi...	health	health	0.139
Los condilomas o verrugas...	health	health	0.139
Una inteligencia artifici...	health	science	0.122
15 recetas con nueces par...	health	health	0.109
A la caza del segundoaguj...	health	science	0.104
Las bacterias intestinale...	health	health	0.096
Detectan envejecimiento p...	health	health	0.091
Sánchez aspira a que la U...	health	politics	0.075

Tabla 6: Documentos clasificados en la categoría *health* ordenados por similaridad con el glosario

Título	Clasificación	Categoría original	Similaridad
Cuándo es la última lunar...	science	science	0.39
Un ordenador cuánticologr...	science	science	0.281
¿Podían los diplodocuspro...	science	science	0.28
Marte no está muerto:hall...	science	science	0.262
Un meteorito sobrevuelala...	science	science	0.244
Medio siglo desde el Apol...	science	science	0.237
¿Por qué se echa sal a la...	science	science	0.223
La catástrofe del vacío (...)	science	science	0.215
PSOE y Unidas Podemos pla...	science	politics	0.214
Los responsables de medir...	science	science	0.2
Lluvia de estrellas de la...	science	science	0.189
Nudos que representan núm...	science	science	0.186
¿Y si los dinosaurios no ...	science	science	0.186
MBST® Resonancia Terapéut...	science	health	0.179
Buenos hábitos para contr...	science	health	0.17
Así es el síndrome de la ...	science	science	0.168
La ciencia confirma que e...	science	health	0.143
Hallada en Marte una zona...	science	science	0.139
Xavier Trias llama a los ...	science	politics	0.139
¿Existen minisatélites qu...	science	science	0.114
El mejor superalimento fi...	science	health	0.107
Stakhovsky carga contra l...	science	sports	0.094

Tabla 7: Documentos clasificados en la categoría *science* ordenados por similaridad con el glosario

Título	Clasificación	Categoría original	Similaridad
El Ministerio de Igualdad...	politics	politics	0.375
El espionaje a periodista...	politics	politics	0.324
La vicepresidenta de Ceut...	politics	politics	0.293
Que los politicos noJalte...	politics	politics	0.282
Los 5 alimentos populares...	politics	health	0.265
El Poder Judicial propone...	politics	politics	0.248
David De la Cruz, un cand...	politics	politics	0.247
Perú teme una debacle del...	politics	politics	0.225
Feijóo rechaza la moción ...	politics	politics	0.204
La comisión de Justicia a...	politics	politics	0.198
Los conservadores del Pod...	politics	politics	0.194
Podemos propone un cheque...	politics	politics	0.194
Sandro Rosell anunciará t...	politics	politics	0.193
Podemos ve con temor la r...	politics	politics	0.191
Que el ruido no ensordezc...	politics	politics	0.183
China, Ucrania y el litio...	politics	science	0.18
España y el bloque de paí...	politics	politics	0.174
No te obsesiones con laNo...	politics	science	0.173
Una senadora demócrata de...	politics	politics	0.172
El Gobierno cambia la ley...	politics	politics	0.172
El PSOE y Unidas Podemos ...	politics	politics	0.154
Leclerc explica por qué s...	politics	sports	0.15
Trias se suma a la carrer...	politics	politics	0.145
Audi crece para la F1 - A...	politics	sports	0.143
EE.UU. anunciará el marte...	politics	science	0.143
Campazzo quiere jugar en ...	politics	sports	0.142
El ADN más antiguo jamás ...	politics	science	0.139
Hito histórico de EEUU en...	politics	science	0.134
UPN rompe con PP y Ciudad...	politics	politics	0.112
En defensa del voto a los...	politics	politics	0.083

Tabla 8: Documentos clasificados en la categoría *politics* ordenados por similaridad con el glosario

Título	Clasificación	Categoría original	Similaridad
Meloni cancela su presenc...	No clasificado	politics	0.0
Álex Palou, confirmado co...	No clasificado	sports	0.0
La respuesta de Rafa Nada...	No clasificado	sports	0.0
La dieta que reduce en un...	No clasificado	health	0.0
Doncic se sale en la Meca...	No clasificado	sports	0.0
El descalabro gigante de ...	No clasificado	sports	0.0
La lesión de Juancho Hern...	No clasificado	sports	0.0
Seis niños mueren en el R...	No clasificado	health	0.0
El H2Med costará 2.500 mi...	No clasificado	politics	0.0

Tabla 9: Documentos no clasificados

## G. Otros clasificadores en detalle

El código fuente que hemos empleado para implementar otro tipo de clasificadores se encuentra en [Clasificacion-Resto.ipynb](#)

Comenzamos cargando el tokenizador a partir de la librería *spacy* y definimos las funciones para crear las diferentes mutaciones de datos.

```
1 nlp = spacy.load("es_core_news_lg") #Mejor modelo optimizado para la CPU
2
3 #BoW
4 bow_vectorizer = CountVectorizer(min_df=0.01, max_df=0.9, vocabulary=vocabulario_train)
5
6
7
8 #Tf-idf
9 tfidf_vectorizer = TfidfTransformer()
10
11 #Funciones de WV.
12 def averaged_word_vectorizer(corpus):
13     '''Aplica la función de cálculo del WE promedio a todos los
14     documentos del corpus (cada doc es una lista de tokens)'''
15     features = [nlp(doc).vector
16                 for doc in corpus]
17     return np.array(features)
18
19 def tfidf_wtd_avg_word_vectors(doc, word_tfidf_map):
20     '''Aplica la función de cálculo del WE ponderado por TF-IDF
21     a un documento (como lista de tokens)'''
22     tokens = doc.split()
23
24     feature_vector = np.zeros((nlp.vocab.vectors_length,), dtype="float64")
25     wts = 0.
26     for word in tokens:
27         if nlp.vocab[word].has_vector and word_tfidf_map.get(word, 0): #s lo considera
28             palabras conocidas
29             weighted_word_vector = word_tfidf_map[word] * nlp.vocab[word].vector
30             wts = wts + 1
31             feature_vector = np.add(feature_vector, weighted_word_vector)
32     if wts:
33         feature_vector = np.divide(feature_vector, wts)
34
35     return feature_vector
36
37 def tfidf_weighted_averaged_word_vectorizer(corpus, word_tfidf_map):
38     '''Aplica la función de cálculo del WE ponderado por TF-IDF a todos los
39     documentos del corpus (cada doc es una lista de tokens)'''
40     features = [tfidf_wtd_avg_word_vectors(doc, word_tfidf_map)
41                 for doc in corpus]
42     return np.array(features)
```

Código 12: Mutaciones de los datos

De esta forma, quedaría:

```
1 # caracter sticas bag of words
2 bow_train_features = bow_vectorizer.fit_transform(norm_train_corpus)
3 bow_test_features = bow_vectorizer.transform(norm_test_corpus)
4
5 # caracter sticas tfidf (a partir del BoW)
6 tfidf_train_features = tfidf_vectorizer.fit_transform(bow_train_features)
7 tfidf_test_features = tfidf_vectorizer.transform(bow_test_features)
8
9 # caracter sticas averaged word vector
10 avg_wv_train_features = averaged_word_vectorizer(norm_train_corpus)
11 avg_wv_test_features = averaged_word_vectorizer(norm_test_corpus)
12
13 # caracter sticas tfidf weighted averaged word vector
14 word_tfidf_map = {key:value for (key, value) in zip(bow_vectorizer.get_feature_names_out(),
15 , tfidf_vectorizer.idf_)}
```

```

15
16 tfidf_wv_train_features = tfidf_weighted_averaged_word_vectorizer(norm_train_corpus,
17 word_tfidf_map)
18
19 tfidf_wv_test_features = tfidf_weighted_averaged_word_vectorizer(norm_test_corpus,
20 word_tfidf_map)
21

```

Código 13: Diferentes tipos de entradas

Denifimos ahora las diferentes métricas a emplear y la implementación del entrenamiento y evaluación.

```

1 def get_metrics(true_labels, predicted_labels):
2     """Calculamos distintas metricas sobre el
3     rendimiento del modelo. Devuelve un diccionario
4     con los parametros medidos"""
5
6     return {
7         'Accuracy': np.round(
8             metrics.accuracy_score(true_labels,
9                                     predicted_labels),
10            3),
11         'Precision': np.round(
12             metrics.precision_score(true_labels,
13                                     predicted_labels,
14                                     average='weighted',
15                                     zero_division=0),
16            3),
17         'Recall': np.round(
18             metrics.recall_score(true_labels,
19                                   predicted_labels,
20                                   average='weighted',
21                                   zero_division=0),
22            3),
23         'F1 Score': np.round(
24             metrics.f1_score(true_labels,
25                               predicted_labels,
26                               average='weighted',
27                               zero_division=0),
28            3)}
29
30
31 def train_predict_evaluate_model(classifier,
32                                   train_features, train_labels,
33                                   test_features, test_labels):
34     """Funcion que entrena un modelo de clasificacion sobre
35     un conjunto de entrenamiento, lo aplica sobre un conjunto
36     de test y devuelve la prediccion sobre el conjunto de test
37     y las metricas de rendimiento"""
38     # genera modelo
39     classifier.fit(train_features, train_labels)
40     # predice usando el modelo sobre test
41     predictions = classifier.predict(test_features)
42     # eval a rendimiento de la predicci n
43     metricas = get_metrics(true_labels=test_labels,
44                             predicted_labels=predictions)
45     return predictions, metricas
46

```

Código 14: Definición de métricas y proceso

Pasamos a la primera fase donde probamos los distintos tipos de clasificadores del paquete *Scikit-Learn* con los distintos tipos de datos.

```

1 modelLR = LogisticRegression(solver='liblinear')
2 modelNB = GaussianNB() #var_smoothing=1e-9
3 modelSVM = SGDClassifier(loss='hinge', max_iter=1000)
4 modelRBF SVM = SVC(gamma='scale', C=2)
5
6 modelos = [('Logistic Regression', modelLR),
7            ('Naive Bayes', modelNB),
8            ('Linear SVM', modelSVM),
9            ('Gauss kernel SVM', modelRBF SVM)]

```

```

10
11 metricas = []
12 resultados = []
13
14 # Modelos con caracter sticas BoW
15 bow_train_features_ = bow_train_features.toarray()
16 bow_test_features_ = bow_test_features.toarray()
17 for m, clf in modelos:
18     prediccion, metrica = train_predict_evaluate_model(classifier=clf,
19                                                         train_features=bow_train_features_,
20                                                         train_labels=train_labels,
21                                                         test_features=bow_test_features_,
22                                                         test_labels=test_labels)
23     metrica['modelo']=f'{m} BoW'
24     resultados.append(prediccion)
25     metricas.append(metrica)
26
27 # Modelos con caracter sticas TF-IDF
28 tfidf_train_features_ = tfidf_train_features.toarray()
29 tfidf_test_features_ = tfidf_test_features.toarray()
30 for m, clf in modelos:
31     prediccion, metrica = train_predict_evaluate_model(classifier=clf,
32                                                         train_features=tfidf_train_features_,
33                                                         train_labels=train_labels,
34                                                         test_features=tfidf_test_features_,
35                                                         test_labels=test_labels)
36     metrica['modelo']=f'{m} tfidf'
37     resultados.append(prediccion)
38     metricas.append(metrica)
39
40 # Modelos con caracter sticas averaged word vectors
41 avg_wv_train_features_ = avg_wv_train_features
42 avg_wv_test_features_ = avg_wv_test_features
43 for m, clf in modelos:
44     prediccion, metrica = train_predict_evaluate_model(classifier=clf,
45                                                         train_features=avg_wv_train_features_,
46                                                         train_labels=train_labels,
47                                                         test_features=avg_wv_test_features_,
48                                                         test_labels=test_labels)
49     metrica['modelo']=f'{m} averaged'
50     resultados.append(prediccion)
51     metricas.append(metrica)
52
53 # Modelos con caracter sticas tfidf weighted averaged word vectors
54 tfidf_wv_train_features_ = tfidf_wv_train_features
55 tfidf_wv_test_features_ = tfidf_wv_test_features
56 for m, clf in modelos:
57     prediccion, metrica = train_predict_evaluate_model(classifier=clf,
58                                                         train_features=tfidf_wv_train_features_,
59                                                         train_labels=train_labels,
60                                                         test_features=tfidf_wv_test_features_,
61                                                         test_labels=test_labels)
62     metrica['modelo']=f'{m} tfidf wv'
63     resultados.append(prediccion)
64     metricas.append(metrica)
65

```

Código 15: Definición de modelos y entrenamiento con diferentes entradas

Obtenemos los resultados de la primera fase de la siguiente forma:

```

1 metricas = pd.DataFrame(metricas)
2 metricas = metricas.sort_values('Accuracy',ascending=False)
3 display(metricas)
4
5 metricas.head(1)
6

```

Código 16: Resultados de la primera fase

Pamos pues a la segunda fase escogiendo los datos que mejor funcionan y los cuatro mejores modelos. Generalmente, son los promediados de los word embeddings. Entrenamos los modelos escogidos probando diferentes parámetros para comprobar si su funcionamiento mejora.

```

1  modelos = metricas['modelo'].tolist()
2
3  mejor = modelos[0]
4  segundo = modelos[1]
5  tercero = modelos[2]
6  cuarto = modelos[3]
7
8  moa = [mejor, segundo, tercero, cuarto]
9  print(moa)
10
11

```

Código 17: Selección de los mejores modelos

Como los modelos tienen en el nombre los datos con los que han sido entrenados, es necesario separar dicho nombre para obtener por un lado el nombre de modelo y por otro el conjunto de datos.

```

1  def separar_modelo(modelo):
2      """ funcion para separar del nombre de modelo los datos con los que ha entrenado
3      - modelo: nombre del modelo
4      """
5      sep = modelo.split(' ')
6
7      if sep[len(sep)-1] == "wv":
8          mo = ' '.join(sep[:len(sep)-2])
9
10     else:
11         mo = ' '.join(sep[:len(sep)-1])
12
13
14     datos = sep[len(sep)-1]
15
16
17     return mo, datos
18
19
20
21
22     dat = []
23     mod = []
24
25     for i in moa:
26         m, d = separar_modelo(i)
27         mod.append(m)
28         dat.append(d)
29
30
31     print(mod)
32     print(dat)
33

```

Código 18: Selección de los mejores modelos

Implementación de la segunda fase:

```

1  metricas2 = []
2  resultados = []
3
4
5  for datos, mo in zip(dat, mod):
6      # DATOS
7      if datos == 'Bow':
8          train_features_ = bow_train_features.toarray()
9          test_features_ = bow_test_features.toarray()
10
11     if datos == 'tfidf':
12         train_features_ = tfidf_train_features.toarray()
13         test_features_ = tfidf_test_features.toarray()
14
15     if datos == 'averaged':
16         train_features_ = avg_wv_train_features
17         test_features_ = avg_wv_test_features
18

```

```

19     else: # tfidf wv
20         train_features_ = tfidf_wv_train_features
21         test_features_ = tfidf_wv_test_features
22
23
24
25
26     # MODELOS
27
28     if mo == 'Logistic Regression':
29         # https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.
LogisticRegression.html
30
31         sol = ["lbfgs", "liblinear", "newton-cg", "newton-cholesky", "sag", "saga"]
32
33         for i in sol:
34             modelLR = LogisticRegression(solver=i)
35             prediccion, metrica = train_predict_evaluate_model(classifier=modelLR,
36                                                                 train_features=train_features_,
37                                                                 train_labels=train_labels,
38                                                                 test_features=test_features_,
39                                                                 test_labels=test_labels)
40             metrica['modelo']=f'{i} - Logistic Regression ' + datos
41             resultados.append(prediccion)
42             metricas2.append(metrica)
43
44         if mo == 'Naive Bayes':
45             # https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB
.html
46             var_smo = [1e-1,1e-2,1e-3,1e-4,1e-5,1e-6,1e-7,1e-8,1e-9,1e-10,1e-11,1e-12,1e-13,1e
-14,1e-15]
47
48             for i in var_smo:
49                 modelNB = GaussianNB(var_smoothing=i)
50                 prediccion, metrica = train_predict_evaluate_model(classifier=modelNB,
51                                                                 train_features=train_features_,
52                                                                 train_labels=train_labels,
53                                                                 test_features=test_features_,
54                                                                 test_labels=test_labels)
55                 metrica['modelo']=f'{i} - Naive Bayes ' + datos
56                 resultados.append(prediccion)
57                 metricas2.append(metrica)
58
59             if mo == 'Linear SVM':
60                 # https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.
SGDClassifier.html
61                 loss = ["hinge", "log_loss", "log", "modified_huber", "squared_hinge", "perceptron
", "squared_error", "huber", "epsilon_insensitive", "squared_epsilon_insensitive"]
62
63                 for i in loss:
64                     modelSVM = SGDClassifier(loss=i, max_iter=1000)
65                     prediccion, metrica = train_predict_evaluate_model(classifier=modelSVM,
66                                                                 train_features=train_features_,
67                                                                 train_labels=train_labels,
68                                                                 test_features=test_features_,
69                                                                 test_labels=test_labels)
70                     metrica['modelo']=f'{i} - Linear SVM ' + datos
71                     resultados.append(prediccion)
72                     metricas2.append(metrica)
73
74
75             if mo == 'Gauss kernel SVM':
76                 # https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html
77                 gamma = ["scale", "auto"]
78
79                 for i in gamma:
80                     modelRBFSVM = SVC(gamma=i, C=2)
81                     prediccion, metrica = train_predict_evaluate_model(classifier=modelRBFSVM,
82                                                                 train_features=train_features_,
83                                                                 train_labels=train_labels,
84                                                                 test_features=test_features_,

```

```

85                                     test_labels=test_labels)
86     metrica['modelo']=f'{i} - Gauss kernel SVM ' + datos
87     resultados.append(prediccion)
88     metricas2.append(metrica)
89

```

Código 19: Variación de parámetros para comprobar mejora

Obtenemos los resultados de la segunda fase:

```

1     metricas3 = pd.DataFrame(metricas2)
2     metricas3 = metricas3.sort_values('Accuracy',ascending=False)
3     display(metricas3)
4     metricas3.head(1)
5

```

Código 20: Obtención de resultados de la segunda fase

Pamos a la tercera fase. Por el momento tenemos los modelos que mejor funcionan con el conjunto de datos que mejor se adapta a los modelos. Ya solo nos falta comprobar el rendimiento de estos modelos en la práctica. Seleccionamos los modelos:

```

1     modelos = metricas3['modelo'].tolist()
2
3     mejor = modelos[0]
4     segundo = modelos[1]
5     tercero = modelos[2]
6     # cuarto = modelos[3]
7
8     moa2 = [mejor, segundo, tercero] #, cuarto
9     print(moa2)
10
11
12     dat = []
13     mod = []
14     params = []
15
16     for i in moa2:
17         m, d = separar_modelo(i)
18
19         dat.append(d)
20
21
22         p = m.split('-')
23
24         if 'newton' in p:
25             para = '-'.join(p[0:2])
26             mdo = ' '.join(p[2:])
27
28         else:
29             para = p[0]
30             mdo = ' '.join(p[1:])
31
32
33         params.append(para.strip())
34         mod.append(mdo.strip())
35
36
37     print(mod)
38     print(dat)
39     print(params)
40

```

Código 21: Selección de modelos para la tercera fase

Definimos una función para la puesta en producción de dichos modelos para que no haya que copiar y pegar los modelos sino que lo haga de forma automática.

```

1     def produccion(mod, dat, params):
2         """ Funcion para crear los tres mejores modelos en produccion
3         - mod: lista con los nombres de modelos
4         - dat: lista con los datos de los mejores modelos
5         - params: lista de los parametros correspondientes a los modelos

```



```

6      """
7
8      resultados2 = []
9      metricas4 = []
10     clasificadores = []
11
12     for datos, mo, param in zip(dat, mod, params):
13
14         # DATOS
15         if datos == 'Bow':
16             train_features_ = bow_train_features.toarray()
17             test_features_ = bow_test_features.toarray()
18
19         if datos == 'tfidf':
20             train_features_ = tfidf_train_features.toarray()
21             test_features_ = tfidf_test_features.toarray()
22
23         if datos == 'averaged':
24             train_features_ = avg_wv_train_features
25             test_features_ = avg_wv_test_features
26
27         else: # tfidf wv
28             train_features_ = tfidf_wv_train_features
29             test_features_ = tfidf_wv_test_features
30
31
32
33
34         # MODELOS
35
36         if mo == 'Logistic Regression':
37             # https://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.
38             LogisticRegression.html
39
40             modelLR = LogisticRegression(solver=param)
41             prediccion, metrica = train_predict_evaluate_model(classifier=modelLR,
42                                                                 train_features=train_features_,
43                                                                 train_labels=train_labels,
44                                                                 test_features=test_features_,
45                                                                 test_labels=test_labels)
46
47             metrica['modelo']=f'{param} - Logistic Regression'
48             resultados2.append(prediccion)
49             metricas4.append(metrica)
50             clasificadores.append(modelLR)
51
52         if mo == 'Naive Bayes':
53             # https://scikit-learn.org/stable/modules/generated/sklearn.naive\_bayes.
54             GaussianNB.html
55
56             modelNB = GaussianNB(var_smoothing=param)
57             prediccion, metrica = train_predict_evaluate_model(classifier=modelNB,
58                                                                 train_features=train_features_,
59                                                                 train_labels=train_labels,
60                                                                 test_features=test_features_,
61                                                                 test_labels=test_labels)
62
63             metrica['modelo']=f'{param} - Naive Bayes'
64             resultados2.append(prediccion)
65             metricas4.append(metrica)
66             clasificadores.append(modelNB)
67
68         if mo == 'Linear SVM':
69             # https://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.
70             SGDClassifier.html
71
72             modelSVM = SGDClassifier(loss=param, max_iter=1000)
73             prediccion, metrica = train_predict_evaluate_model(classifier=modelSVM,
74                                                                 train_features=train_features_,
75                                                                 train_labels=train_labels,
76                                                                 test_features=test_features_,
77                                                                 test_labels=test_labels)

```

```

74     metrica['modelo']=f'{param} - Linear SVM'
75     resultados2.append(prediccion)
76     metricas4.append(metrica)
77     clasificadores.append(modelSVM)
78
79
80     if mo == 'Gauss kernel SVM':
81         # https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html
82
83         modelRBFSVM = SVC(gamma=param, C=2)
84         prediccion, metrica = train_predict_evaluate_model(classifier=modelRBFSVM,
85                                                         train_features=train_features_,
86                                                         train_labels=train_labels,
87                                                         test_features=test_features_,
88                                                         test_labels=test_labels)
89
90         metrica['modelo']=f'{param} - Gauss kernel SVM'
91         resultados2.append(prediccion)
92         metricas4.append(metrica)
93         clasificadores.append(modelRBFSVM)
94
95     return resultados2, metricas4, clasificadores
96
97 resultados2, metricas4, clasificadores = produccion(mod, dat, params)
98
99

```

Código 22: Definición de los modelos en una tercera fase

Obtenemos los resultados la tercera fase y cargamos el registro de las anteriores iteraciones para almacenar los resultados de la iteración actual.

```

1  metricas5 = pd.DataFrame(metricas4)
2  mra = pd.read_csv(prediccion_r+'/metricas_produccion.csv')
3
4  tot = pd.merge(mra, metricas5, how='outer')
5
6  tot1 = tot.sort_values('Accuracy',ascending=False)
7  display(tot1)
8
9  predi = []
10 vf = []
11 labs = []
12
13 for c, i in zip(moa2,resultados2):
14     q = [c for n in range(len(i))]
15
16     for f, g, r in zip(i,q, test_labels):
17         predi.append(g)
18         vf.append(f)
19         labs.append(r)
20
21
22 print(len(predi))
23 print(len(vf))
24 print(len(labs))
25
26
27 res = pd.DataFrame()
28 res['modelo'] = vf
29 res['prediccion'] = predi
30 res['category'] = labs
31
32 display(res)
33

```

Código 23: Obtención de resultados de la tercera fase

Por último comprobamos el rendimiento de los tres mejores modelos para cada una de las categorías a partir de la matriz de confusión.

```

1  for title, classifier in zip(moa2, clasificadores):

```

```

2     disp = ConfusionMatrixDisplay.from_estimator(
3         classifier,
4         test_features_,
5         test_labels,
6         display_labels=categorias,
7         cmap=plt.cm.Blues,
8         normalize=None,
9     )
10
11     disp.ax_.set_title(title)
12

```

Código 24: Matriz de confusión de los tres mejores modelos

## Referencias

- [1] Christopher D. Manning, Prabhakar Raghavan e Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, 2008. DOI: [10.1017/CB09780511809071](https://doi.org/10.1017/CB09780511809071).
- [2] Anand V. Saurkar, Kedar G. Pathare y Shweta A. Gode. “An Overview on Web Scraping Techniques and Tools”. En: 2018.