



Universidad
Internacional
Menéndez Pelayo

Máster Universitario en Investigación en Inteligencia Artificial

Curso 2020-2021

**Recuperación y extracción de información,
grafos y redes sociales**

Práctica Bloque II: Recuperación de información y minería de texto

19 de abril de 2021

Laura Rodríguez Navas
DNI: 43630508Z

e-mail: rodrigueznabas@posgrado.uimp.es

Índice

1. Resumen	3
2. Rastreador web (crawler)	3
3. K-Means	5
3.1. Datos de entrada	5
3.2. Palabras vacías, derivación y tokenización	5
3.3. Extracción de características	6
3.4. Entrenamiento del algoritmo	7
3.5. Evaluación	7
Bibliografía	8

1. Resumen

En esta práctica se ha implementado un rastreador web (crawler) en Python [1] (ver Sección 2), que se complementa con un proceso de agrupamiento (ver Sección 3), también implementado en Python, de la información extraída por el rastreador web.

2. Rastreador web (crawler)

En esta sección se describe como se ha implementado el rastreador web (crawler) en Python usando la librería [Scrapy](#). Para empezar con la implementación se debe ejecutar el siguiente comando:

```
$ scrapy startproject books
```

Este comando crea un proyecto Scrapy en el directorio books, siguiendo la [estructura por defecto](#), común para todo proyecto Scrapy. También crea el fichero `scrapy.cfg`, que contiene el nombre del módulo en Python que define la configuración del proyecto books (`books.settings`). El proyecto lo he nombrado books, porqué se rastreará y se recuperará información de un catálogo de libros que se encuentra en la página web: <http://books.toscrape.com>.

Una vez se ha creado el proyecto, se tienen que definir los ítems de cada libro que se quieran extraer del catálogo. En este caso los ítems que se van a extraer son: el título, la categoría, la descripción, el precio y la valoración de cada libro. Para ello, se tiene que modificar el fichero `books/items.py`, para incluir los cinco ítems que se quieren extraer. Vemos el contenido de `items.py` a continuación:

```
import scrapy

class BooksItem(scrapy.Item):
    # define the fields for your item here like:
    # name = scrapy.Field()
    title = scrapy.Field()
    category = scrapy.Field()
    description = scrapy.Field()
    price = scrapy.Field()
    rating = scrapy.Field()
```

El siguiente paso es describir la manera de extraer la información definida en el fichero `items.py`. Para ello, se utilizan reglas de expresión [XPath](#) y [CSS](#). Por ejemplo, si nos fijamos en el código HTML de uno de los libros que se van a rastrear (ver Figura 1), veremos que el título del libro es fácil de extraer con la siguiente regla de expresión CSS: `"h1 ::text"`. Pero cuando la extracción de información se complica un poco más, se usan reglas de expresión XPath. Por ejemplo, para extraer las descripciones de todos los libros se usará la regla de expresión: `"//div[@id='product_description']/following-sibling::p/text()"`. Una vez, definidas todas las reglas de expresión para cada ítem que se va a rastrear, se crea la araña `books/spiders/books_toscrape.py`.

Definición 1. Las arañas son clases que definen cómo se rastrea una página web determinada (o un grupo de páginas web), incluido cómo realizar el rastreo y cómo extraer la información deseada. En otras palabras, las arañas son el lugar donde se define el comportamiento personalizado para rastrear y analizar las páginas web.

En el caso de la práctica, la araña `books.toscrape` representará el lugar donde se definen las reglas de expresión. En las arañas también se tienen que especificar las solicitudes iniciales para rastrear las URLs y una función de devolución de llamada (*parse*) a la que se llamará para generar los ítems de respuesta de esas solicitudes. Por último, los ítems devueltos por las arañas normalmente se conservan en una base de datos o se escriben en un archivo. En el caso de la práctica, los ítems (título, categoría, descripción, precio y valoración de cada libro)

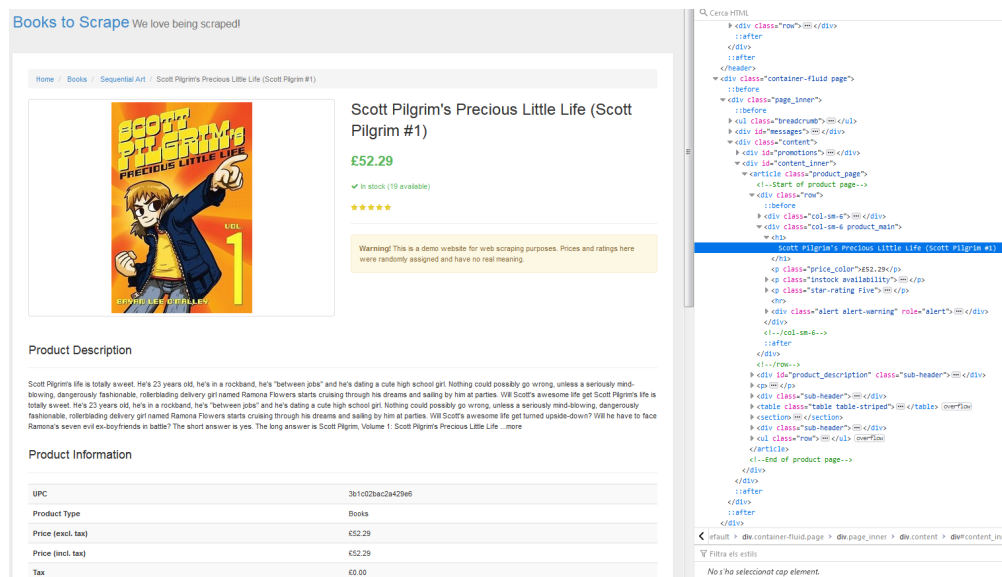


Figura 1: Ejemplo de libro a rastrear.

devueltos por la araña serán guardados en el fichero *books.json*. La araña *books.toscrape* que procesa todas las URLs descubiertas de <http://books.toscrape.com>, utilizando la función *parse* que a su vez llama a la función *parse_book_page* donde son definidas todas las reglas de expresión de cómo extraer la información deseada, se muestra a continuación:

```
import scrapy

class BooksToscrapeSpider(scrapy.Spider):
    name = 'books.toscrape'
    allowed_domains = ['books.toscrape.com']
    start_urls = ['http://books.toscrape.com/']

    def parse(self, response):
        for book_url in response.css("article.product_pod > h3 > a ::attr(href)").extract():
            yield scrapy.Request(response.urljoin(book_url), callback=self.parse_book_page)
        next_page = response.css("li.next > a ::attr(href)").extract_first()
        if next_page:
            yield scrapy.Request(response.urljoin(next_page), callback=self.parse)

    @staticmethod
    def parse_book_page(response):
        item = {}
        product = response.css("div.product_main")
        item['title'] = product.css("h1 ::text").extract_first()
        item['category'] = response.xpath("//ul[@class='breadcrumb']/li[@class='active']/preceding-sibling::li[1]/a/text()").extract_first()
        item['description'] = response.xpath("//div[@id='product_description']/following-sibling::p/text()").extract_first()
        price = response.xpath('//th[text()="Price (incl. tax)"]/following-sibling::td/text()').extract_first()
        item['price'] = price.replace('£', '')
        rating = response.xpath('//*[contains(@class, "star-rating")/@class').extract_first()
        item['rating'] = rating.replace('star-rating ', '')
        yield item
```

En este momento ya podemos iniciar la araña para que recupere la información del catálogo web y la guarde en el fichero *books.json*, aunque primero es recomendable modificar el fichero *books/settings.py* para limitar el acceso de la araña al catálogo web, ya que podemos generar un ataque **DDoS**. Para ello, debemos descomentar la variable **DOWNLOAD_DELAY** y darle un valor en segundos (p.ej. **DOWNLOAD_DELAY = 3**).

Para iniciar la araña se debe ejecutar el siguiente comando:

```
$ cd books
$ scrapy crawl books.toscrape -o books.json
```

3. K-Means

En esta sección se describe como se ha implementado el proceso de agrupamiento en Python (ver directorio `kmeans` en [1]) usando la librería `scikit-learn` [2]. Existen muchos algoritmos de agrupación, y para esta práctica se ha elegido el algoritmo **K-Means**. Concretamente, el algoritmo K-Means agrupará los títulos de los libros del catálogo web (recuperados en la sección 2) en diferentes clústeres por sus similitudes.

3.1. Datos de entrada

Si nos fijamos en `main` del fichero `kmeans/kmeans.py`, vemos que empieza extrayendo la información de los libros almacenada en el fichero `books/books.json` (1000 documentos) y la convierte en un *DataFrame* (ver Definición 2). A continuación, se eliminan los valores NaN que pudieran existir en él y es almacenado en un fichero CSV (`kmeans/books.csv`). Para ello, se usa la librería `pandas` [3].

Las primeras líneas del *DataFrame* que forman los datos de entrada:

title	category	description	price	rating
Sapiens: A Brief History of Humankind	History	From a renowned historian...	54.23	Five
Sharp Objects	Mystery	WICKED above her hipbone, GIRL...	47.82	Four
Soumission	Fiction	Dans une France assez...	50.10	One
Tipping the Velvet	Historical Fiction	Erotic and absorbing...	53.74	One
A Light in the Attic	Poetry	It's hard to imagine...	51.77	Three

Tabla 1: Conjunto de datos de entrada.

En este punto, concretamos la estructura de datos que se utilizará para alimentar el algoritmo K-Means. Se crea una lista que contiene los títulos de los libros.

```
titles = df["title"].to_list()
print(titles[:10]) # first 10 titles
>> ['Sapiens: A Brief History of Humankind', 'Sharp Objects', 'Soumission', 'Tipping the
    Velvet', 'A Light in the Attic', "It's Only the Himalayas", 'Libertarianism for
    Beginners', 'Mesaerion: The Best Science Fiction Stories 1800-1849', 'Olio', 'Our Band
    Could Be Your Life: Scenes from the American Indie Underground, 1981-1991']
```

Definición 2. Un *DataFrame* es una estructura de datos bidimensional etiquetada que acepta diferentes tipos de datos de entrada organizados en columnas. Se puede pensar en un *DataFrame* como una hoja de cálculo o una tabla SQL.

3.2. Palabras vacías, derivación y tokenización

Esta sección se centra en definir algunas funciones para manipular los títulos. Primero, se obtiene la lista de palabras vacías en inglés (ver Definición 3), usando la librería `NLTK` [4]. La lista de palabras vacías es en inglés

ya que los títulos de los libros están en inglés. Segundo, se obtiene el *Snowball Stemmer*, también usando la librería NLTK, para descomponer las palabras que componen cada título en su raíz.

```
# nltk's English stopwords as variable called 'stopwords'
stopwords = nltk.corpus.stopwords.words('english')
# nltk's SnowballStemmer as variable called 'stemmer'
stemmer = SnowballStemmer("english")
print(stopwords[:10]) # first 10 stopwords
>> ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're"]
```

Definición 3. Las palabras vacías son palabras sin significado como artículos, pronombres, preposiciones, etc. que son filtradas antes o después del procesamiento de datos en lenguaje natural (texto).

A continuación se define la función *tokenize_and_stem*:

```
def tokenize_and_stem(text):
    # first tokenize by sentence, then by word to ensure that punctuation is caught as it's
    # own token
    tokens = [word for sent in nltk.sent_tokenize(text) for word in nltk.word_tokenize(sent)]
    filtered_tokens = []
    # filter out any tokens not containing letters (e.g., numeric tokens, raw punctuation)
    for token in tokens:
        if re.search('[a-zA-Z]', token):
            filtered_tokens.append(token)
    stems = [stemmer.stem(t) for t in filtered_tokens]
    return stems
```

Esta función, que divide los títulos en una lista de tokens y los descompone con su raíz, será usada en la siguiente sección para calcular la matriz tf-idf.

3.3. Extracción de características

En esta sección, se calcula la matriz **tf-idf**. Para ello, primero se tienen que contar las apariciones de las palabras por cada documento, que se transforman después en la matriz de documentos y términos (ver Figura 2). Después, las palabras que aparecen con frecuencia dentro de un documento pero no con frecuencia dentro del **corpus** reciben una ponderación más alta, ya que se supone que estas palabras contienen más significado en relación con el documento.

Definición 4. En lingüística, un *corpus* o *corpus de texto* es un recurso lingüístico formado por un conjunto grande de textos y estructurados.

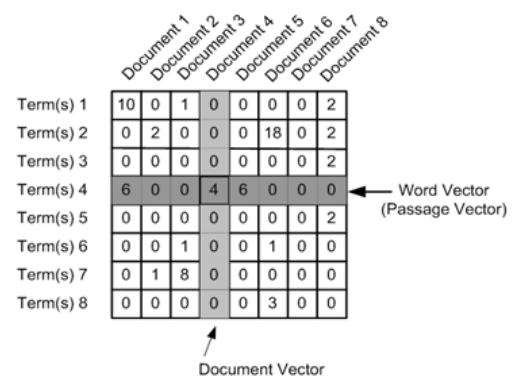


Figura 2: Matriz de documentos y términos.

La matriz tf-idf para todos los títulos se calcula usando la librería `sklearn.feature_extraction.text.TfidfVectorizer`.

Un par de cosas a tener en cuenta sobre los parámetros definidos en la función `TfidfVectorizer` son:

- `max_features`: construye un vocabulario que solo considera las características máximas principales ordenadas por frecuencia de términos de todo el corpus (ver Definición 4).

- `use_idf`: habilita la reponderación de frecuencia de documentos inversa.
- `gram_range`: define el límite inferior y superior del rango de n-valores para diferentes n-gramos que se extraerán. Principalmente se usa para observar unigramas, bigramas y trigramas. Ver [n-grams](#).

```
# define vectorizer parameters
tfidf_vectorizer = TfidfVectorizer(max_features=200000, stop_words=stopwords, use_idf=True,
    tokenizer=tokenize_and_stem, gram_range=(1, 3))
tfidf_matrix = tfidf_vectorizer.fit_transform(titles) # fit the vectorizer to titles
print(tfidf_matrix.shape)
>> (998, 7258)
```

Vemos que la matriz tf-idf está formada por 998 documentos y 7258 términos.

3.4. Entrenamiento del algoritmo

Ahora pasemos a la parte divertida. Usando la matriz tf-idf, calculada en la sección anterior, se entrena el algoritmo K-Means utilizando [sklearn.cluster.KMeans](#). K-means se inicializa con un número predeterminado de clústeres. Elegí 40 como número predeterminado de clústeres, ya que el conjunto de datos contiene libros que pertenecen a una de las 40 categorías, descartando 10 categorías que solo hacen referencia a un único libro.

```
# K-Means clustering
num_clusters = 40
km = KMeans(n_clusters=num_clusters)
km.fit(tfidf_matrix)
clusters = km.labels_.tolist()
```

Para analizar los resultados del algoritmo, primero creamos un nuevo *Dataframe* con los títulos y su clústeres asignados. Vemos las primeras líneas del contenido del nuevo *Dataframe* a continuación:

title	cluster
Sapiens: A Brief History of Humankind	26
Sharp Objects	0
Soumission	12
Tipping the Velvet	12
A Light in the Attic	9

Tabla 2: Dataframe de los títulos y su clústeres asignados.

3.5. Evaluación

Para la visualización de los clústeres creados en la sección anterior, se tienen que hacer algún procesamiento.

Bibliografía

- [1] Laura Rodríguez-Navas. Recuperación de información y minería de texto. <https://github.com/lrodrin/masterAI/tree/master/A14>, 2021.
- [2] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [3] Jeff Reback, Wes McKinney, jbrockmendel, Joris Van den Bossche, Tom Augspurger, Phillip Cloud, gyoung, Sinhrks, Simon Hawkins, Matthew Roeschke, Adam Klein, Terji Petersen, Jeff Tratner, Chang She, William Ayd, Shahar Naveh, Marc Garcia, Jeremy Schendel, Andy Hayden, Daniel Saxton, Vytutas Jancauskas, Ali McMaster, Pietro Battiston, Skipper Seabold, patrick, Kaiqi Dong, chris b1, h vetinari, Stephan Hoyer, and Marco Gorelli. pandas-dev/pandas: Pandas 1.1.5, December 2020.
- [4] Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing text with the natural language toolkit*. "O'Reilly Media, Inc.", 2009.