

Scraping Fitness

Estudiante: Jaime Cabero Creus

Estudiante: David García Gondell

A Coruña, December de 2021.

Índice general

1	Introducción	2
1.1	Contextualización y descripción del sistema	2
2	Tecnologías empleadas	3
2.1	Servicio	3
2.1.1	Python	3
2.1.2	Scrapy	3
2.1.3	Elasticsearch	4
2.2	Aplicación web SPA	4
2.2.1	React	4
2.2.2	Redux	4
2.2.3	Material UI	4
3	Desarrollo	5
3.1	Funcionalidades	5
3.2	Servicio	5
3.2.1	Fichero de configuración	6
3.2.2	Fichero <i>DiaSpider</i>	6
3.2.3	Fichero <i>pipeline</i>	7
3.3	Interfaz de usuario	8

Introducción

ESTE primer capítulo de la memoria incluirá una explicación del sistema desarrollado y de los objetivos que buscamos conseguir con este proyecto.

1.1 Contextualización y descripción del sistema

Para este proyecto se han considerado multitud de dominios sobre los que extraer datos, para ello hemos tenido en cuenta que tuvieran una cantidad de información suficiente, que esta fuera accesible y por último que el propio dominio no contase con suficientes mecanismos para tratar esos datos. Finalmente, hemos decidido por extraer datos de una cadena de supermercados para proporcionar diferentes dietas y búsquedas de alimentos a nuestros usuarios.

Una vez elegido el dominio, se ha pensado en como hacer el proceso de extracción y tratamiento de los datos. El sistema ideado constará de un scraper en lenguaje Python que se encarga de recopilar información de la página web de supermercados Dia, de un servicio Elasticsearch para indexar la información y de una aplicación web SPA construida utilizando React que se encarga de mostrar las búsquedas en Elasticsearch.

Tecnologías empleadas

EN este capítulo se expondrán las tecnologías empleadas en este proyecto. Se explicará el motivo de su elección y el uso dado a las mismas en el proceso de desarrollo.

2.1 Servicio

2.1.1 Python

Python es un lenguaje de programación de alto nivel. Se trata de un lenguaje multiplataforma y de código abierto, que no necesita de un compilador que lo traduzca a lenguaje máquina, sino que se sirve de un interpretador para ejecutarse. En este proyecto se ha empleado Python para la implementar el servicio ya que dispone del framework Scrapy para recolectar los datos que posteriormente vamos a tratar.

2.1.2 Scrapy

Scrapy es un framework de *scraping* y *crawling* de código abierto escrito en lenguaje Python; en nuestro proyecto lo utilizamos para realizar *web scraping*, que es una técnica de extracción de información de páginas web de manera automatizada. Scrapy permite declarar clases "Spider" de forma sencilla que se encargan de obtener la información de un tipo de dato disponible en un dominio o conjunto de dominios. Además, una vez obtenidos los datos se puede configurar un pipeline para tratar todos los documentos de uno en uno y enviarlos a otra plataforma, en nuestro caso Elasticsearch.

Decidimos emplear Scrapy frente a otras alternativas porque recibimos una charla sobre como utilizarlo y nos pareció que su funcionamiento era sencillo; además hay que destacar que la documentación existente de Scrapy es rica y está bien estructurada.

2.1.3 Elasticsearch

Elasticsearch es un motor de búsquedas de código abierto basado en Lucene. Este motor cuenta con APIs que nos permiten crear índices de búsqueda desde el servicio y, una vez contruidos, realizar consultas de todo tipo en formato JSON sobre estos índices. Decidimos emplear Elasticsearch debido a que nos pareció sencilla la API que exponía y cómodo su sistema de consultas. Además su API esta disponible tanto para python como para javascript, lo que facilita mucho su uso.

2.2 Aplicación web SPA

2.2.1 React

React es una biblioteca para la construcción de interfaces de usuario en el lenguaje de programación *JavaScript*. Está basado en componentes encapsulados que gestionan su propio estado y tiene un enfoque declarativo. Funciona actualizando y renderizando los componentes correctos de manera eficiente solo cuando los datos cambian. Decidimos emplear React en este proyecto debido a que es una tecnología familiar al equipo de desarrollo y nos permite crear interfaces de usuario de manera rápida y sencilla.

2.2.2 Redux

Redux es una herramienta creada para la gestión del estado en aplicaciones programadas en *JavaScript*, además funciona de manera óptima en conjunto con React, otra de las tecnologías empleadas en este proyecto. Se tomó la decisión de utilizar Redux en este proyecto debido a que al emplear React para el desarrollo de la aplicación web SPA, es muy recomendable realizar una gestión activa del estado de los componentes, para garantizar un comportamiento eficiente de la aplicación.

2.2.3 Material UI

Material UI es una librería de componentes de React con una estética común y una utilidad determinada, que están listos para ser integrados en cualquier proyecto que emplee React. Su diseño está basado en *Material Design*, una serie de especificaciones y guías publicadas por Google para la construcción de interfaces de usuario en aplicaciones Web y móviles. Decidimos utilizar Material UI para darle un aspecto visual atractivo y homogéneo a la aplicación. Además, el equipo de desarrollo está familiarizado con esta tecnología, con lo que se consigue una interfaz de usuario estética en tiempo récord.

Capítulo 3

Desarrollo

EN este apartado de la memoria se describirán las funcionalidades y el proceso de desarrollo del proyecto. En él, se abarcará el desarrollo del servicio, el proceso de *scraping* e indexación de los resultados obtenidos y por último el tratamiento de los datos indexados y su representación en una interfaz de usuario con cierta funcionalidad disponible.

3.1 Funcionalidades

Las funcionalidades desarrolladas en esta práctica son las siguientes:

- Búsqueda de productos por palabras clave
- Búsqueda de productos por precio
- Búsqueda de productos por información nutricional
- Comprar producto (lleva a la página correspondiente en la web de Dia)
- Recomendación de productos de dieta
- Recomendación de productos para dieta de musculación

Estas funcionalidades se explican más en profundidad a continuación.

3.2 Servicio

Vamos a comenzar hablando del desarrollo del servicio. En primer lugar, se ha creado un proyecto de Python con el comando **"scrapy startproject riwsScraper"**.

Este comando creó un directorio "riwsScraper" con el siguiente contenido:

```
1 riwsScraper/  
2   scrapy.cfg          # deploy configuration file  
3  
4   riwsScraper/        # project's Python module, you'll  
   import your code from here  
5       __init__.py  
6  
7       items.py        # project items definition file  
8  
9       middlewares.py  # project middlewares file  
10  
11      pipelines.py     # project pipelines file  
12  
13      settings.py      # project settings file  
14  
15      spiders/         # a directory where you'll later put your  
   spiders  
16          __init__.py
```

Dentro de este directorio se pueden ver diversas clases en Python que ha creado Scrapy. Estas se reparten entre clases de configuración, la clase que sirve de pipeline, y un directorio que contendrá las clases *Spider* que son las encargadas de realizar el proceso de *crawling*, como ya se explicó anteriormente.

3.2.1 Fichero de configuración

El fichero *settings.py* es el encargado de albergar la configuración del servicio. En él viene una configuración inicial que se puede modificar según las necesidades del usuario.

En nuestro caso le pusimos un nombre a nuestro bot, le indicamos donde se ubicaba nuestro módulo de *Spiders* y configuramos nuestro pipeline. Además, limitamos el número de peticiones concurrentes a 8 y le introducimos un pequeño *delay* de un segundo. También, marcamos que nuestro bot siguiera las indicaciones de los ficheros *robot.txt* de las páginas web, los cuales indican a los bots qué URLs del sitio web son las que tienen permitido acceder.

Por último, nos identificamos con el *user-agent*: "Mozilla/5.0(Macintosh; IntelMacOSX10.12.6) AppleWebKit/537.36(KHTML, likeGecko)Chrome/61.0.3163.100Safari/537.36" Todo esto lo hemos hecho de manera preventiva para que no nos identificase como un bot de origen desconocido y nos bloqueasen el acceso.

3.2.2 Fichero *DiaSpider*

Es una clase que extiende la lógica de la clase *CrawlSpider*, y a la que se le da un nombre, un dominio y unas direcciones donde empezar a el proceso de *crawling*.

Este proceso sigue una reglas para discriminar lo que se considera un *item* de lo que no. Estas reglas están construidas como expresiones regulares que nos permiten describir la estructura modelo de una URL que se corresponde con el esquema de dirección de lo que se considera un *item*.

Cada uno de estos *items* los *parsea* utilizando una función *parse – item*, donde en nuestro caso asignamos a una serie de variables, el contenido de los elementos HTML de la pagina correspondiente a un *item*. Con estas variables construimos objetos, que devolvemos en formato JSON para poder indexarlos más adelante con la API de Elasticsearch.

3.2.3 Fichero *pipeline*

Este fichero es una clase que se inicializa al mientras se recopila la información de la web creando un indice de Elasticsearch mediante su API de Pyhton, en este caso llamado *products*. Luego, mediante esta API se indexan todos los documentos que encuentra scrapy durante el proceso de *crawling* utilizando un método que recibe cada *item* como parámetro de forma automática. En nuestro caso, cada *item* tiene un ID que corresponde a la URL del producto y como contenido los siguientes elementos:

- Nombre
- Subcategorías del producto
- Precio
- Ausencia de gluten
- KJ
- Kcal
- Grasas
- Grasas saturadas
- Carbohidratos
- Azúcares
- Proteínas
- Sal
- Ratio de gimnasio (Kcal/Proteínas)

3.3 Interfaz de usuario

Para el desarrollo de la interfaz de usuario se ha creado un proyecto de React mediante el comando **create-react-app**. Con este proyecto básico creado, hemos instalado las dependencias necesarias para poder emplear las tecnologías de Redux y Material UI.

Una vez estuvo lista esta configuración inicial, decidimos lo que íbamos a mostrar con los datos obtenidos. Para ello, primero comenzamos estableciendo dos componentes separados en la aplicación, uno dedicado a mostrar los productos obtenidos y el otro dedicado a realizar búsquedas y aplicar filtros sobre estos productos. Además, en los que cuenten con información nutricional, esta podrá consultarse en formato de tabla al presionar el botón del mismo nombre. El resultado puede apreciarse en la [Figura 3.1](#).



Figura 3.1: Ejemplo de producto completo

En cuanto a los filtros que dispone la aplicación, esta se pensó para que contara con tres modos: el normal, el modo de dieta y el modo de musculación. Estos modos tienen en común un filtro por nombre, otro por precio y tres filtros que se consideran básicos, que nos permiten

obtener productos sin gluten para celíacos, productos sin sal (menos de 0.12g), o productos sin azúcar (menos de 0.5g).

El modo normal, es en el que se encuentra por defecto. Contendrá todos los productos obtenidos en el proceso de scraping y permitirá al usuario utilizar todos los filtros de los que dispone la aplicación. Esta consulta se hace con una query booleana en la que se filtra por nombre y subcategorías con la siguiente sentencia *should*:

```
1 should: [  
2   {  
3     match: {  
4       name: {  
5         query: nameSearch,  
6         fuzziness: "AUTO",  
7         zero_terms_query: "all"  
8       }  
9     }  
10  },  
11  {  
12    match_bool_prefix: {  
13      name: {  
14        query: nameSearch  
15      }  
16    }  
17  },  
18  {  
19    match: {  
20      subcategories: {  
21        query: nameSearch,  
22        fuzziness: "AUTO",  
23        zero_terms_query: "all"  
24      }  
25    }  
26  },  
27  {  
28    match_bool_prefix: {  
29      subcategories: {  
30        query: nameSearch  
31      }  
32    }  
33  },  
34 ]
```

Los filtros de información nutricional en cambio, se hacen gracias a la sentencia *filter* de

la siguiente forma:

```

1 filter: [
2   {
3     range: {
4       glutenLess: {
5         gte: noGlutenCheck ? true : false,
6         lte: true
7       }
8     }
9   },
10  .
11  .
12  .
13 ]

```

Además, la vista que se encarga de mostrar este modo se puede ver en la figura [Figura 3.2.](#):

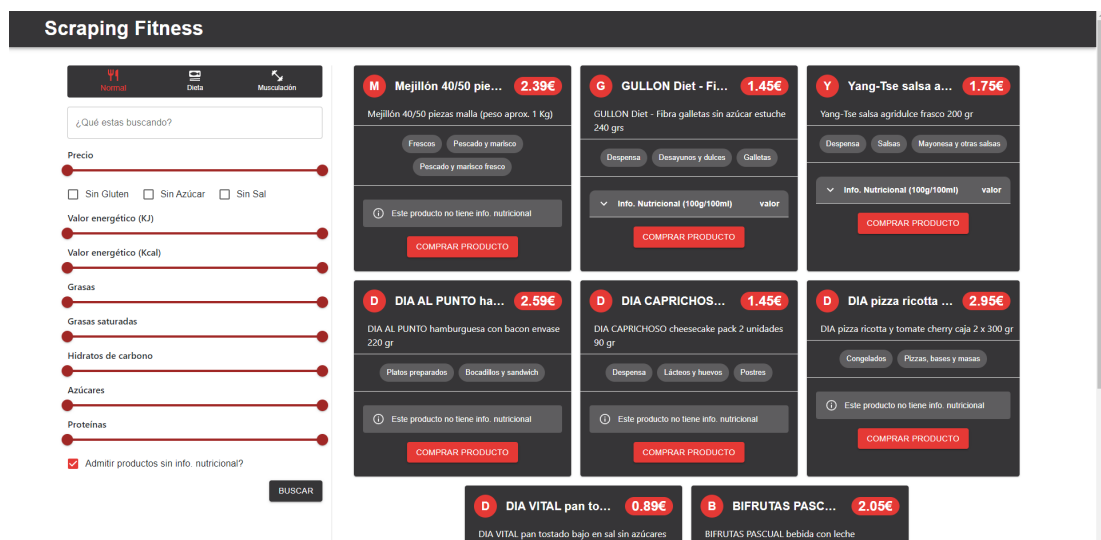


Figura 3.2: Aplicación en modo normal

El modo de dieta, muestra productos que cuenten con un número reducido de calorías, azúcares, grasas y grasas saturadas. En este modo también se podrá filtrar por precio y restringir los alimentos con gluten, sal o azúcar. Estos productos se recuperan de Elasticsearch empleando la siguiente consulta:

Para la consulta del modo dieta se emplea una query function score, que nos permite cambiar el score que le da Elasticsearch a los documentos poniéndole pesos a subconjuntos de los mismos. Un ejemplo de la consulta sería el siguiente:

```
1 query: {  
2   function_score: {  
3     query: { ... },  
4     functions: [  
5       {  
6         filter: {  
7           range: {  
8             kcal: {  
9               gte: 0,  
10              lte: 100  
11            }  
12          }  
13        },  
14        weight: 30  
15      },  
16      {  
17        filter: {  
18          range: {  
19            kcal: {  
20              gte: 100,  
21              lte: 200  
22            }  
23          }  
24        },  
25        weight: 15  
26      }  
27    ]  
28  }  
29 }  
30 }  
31 }
```

En la figura [Figura 3.3](#) se puede ver una consulta en el modo dieta donde se filtra además por productos sin gluten.

Por último, en el modo dedicado a la musculación se muestran principalmente productos que tienen una buena relación entre su aporte calórico y proteico, como se puede apreciar en la sección destacada en color azul en la [Figura 3.4](#). Como sucedía en el modo de dieta, se puede filtrar también por precio, y eliminar de la búsqueda los alimentos con gluten, sal o azúcar. La consulta empleada para recuperar estos productos se hace con una query booleana en la que se filtra por el "gimRatio", que no es más que la relación entre las calorías del producto y sus proteínas, lo que da bastante buenos resultados:

Scraping Fitness

Normal

Dieta

Musculación

Precio

☒ Sin Gluten
 ☐ Sin Azúcar
 ☐ Sin Sal

BUSCAR

D

DR. SCHAR pan ...

2.75€

DR. SCHAR pan de molde con cereales SIN GLUTEN bolsa 300 gr

Despensa

Pan

Pan de molde

Sin gluten

Info. Nutricional (100g/100ml)	valor
Valor energético	1044 KJ
Valor energético (Kcal)	249 kcal
Grasas	6.6 g
Grasas saturadas	0.8 g
Hidratos de carbono	38 g
Azúcar	4.1 g
Proteínas	4.5 g
Sal	0.97 g

COMPRAR PRODUCTO

D

DR. SCHAR pan ...

2.75€

DR. SCHAR pan de molde clásico SIN GLUTEN bolsa 300 gr

Despensa

Pan

Pan de molde

Sin gluten

Info. Nutricional (100g/100ml)	valor
Valor energético	981 KJ
Valor energético (Kcal)	233 kcal
Grasas	3.4 g
Grasas saturadas	0.4 g
Hidratos de carbono	43 g
Azúcar	3.3 g
Proteínas	3.5 g
Sal	0.99 g

COMPRAR PRODUCTO

Figura 3.3: Modo de dieta

```

1 {
2   range: {
3     gymRatio: {
4       gte: 0,
5       lte: 15
6     }
7   }
8 }

```

Scraping Fitness

ψ1

Normal

¿Qué estás buscando?

Precio

☐ Sin Gluten

☐ Sin Azúcar

☐ Sin Sal

BUSCAR

D

DANONE ACTIVI...

2.49€

DANONE ACTIVA bifídus semillas chía y almendras 0% M.G pack 4 unidades 120 gr

Dispensa

Lácteos y huevos

Yogures

Info. Nutricional (100g/100ml)

Valor energético	267 KJ
Valor energético (Kcal)	63 kcal
Grasas	0.8 g
Grasas saturadas	0.1 g
Hidratos de carbono	7.3 g
Azúcar	6.1 g
Proteínas	5 g
Sal	0.15 g

COMPRAR PRODUCTO

B

BROOKLYN TOW...

2.95€

BROOKLYN TOWN hamburguesa de carne de raza L bolsa 2 uds 130 gr

Congelados

Carne y pollo

Sin azúcar

Info. Nutricional (100g/100ml)

Valor energético	1009 KJ
Valor energético (Kcal)	243 kcal
Grasas	19 g
Grasas saturadas	7 g
Hidratos de carbono	0.1 g
Azúcar	0.1 g
Proteínas	18 g
Sal	0.53 g

COMPRAR PRODUCTO

B

BROOKLYN TOW...

2.59€

BROOKLYN TOWN hamburguesas de carne retinto L bolsa 2 uds 130 gr

Congelados

Carne y pollo

Sin azúcar

Info. Nutricional (100g/100ml)

Valor energético	896 KJ
Valor energético (Kcal)	215 kcal
Grasas	15.8 g
Grasas saturadas	7.34 g
Hidratos de carbono	0.1 g
Azúcar	0.1 g
Proteínas	18.3 g
Sal	0.27 g

COMPRAR PRODUCTO

Figura 3.4: Modo de musculación

13