

1. Contexto

Discogs

[Explorar](#)
[Mercado online](#)
[Comunidad](#)

[Iniciar sesión](#)
[Registrarse](#)

Mercado online

Todos los artículos

Artículos que quiero

Compras

Carrito

Ajustes del comprador

Has seleccionado:

Formato: Vinyl

Comprar Discos de vinilo

1 - 25 de 45.275.425 < Prev. Siguiente >

	Ordenar por:	Añadido	Estado	Artista	Título	Sello	Vendedor	Precio	
	Molly Hatchet / Meat Loaf / Ted Nugent - Untitled (LP, Comp)	Sello: DeAgostini, Epic	Cat. n.º: IGDA 1173/174	Estado del soporte: Mint (M) ☹	Condición de la funda: Mint (M)	Ver la página de la edición	de.martin-sergio ★★★★★ 100.0%, 278 valoraciones Enviado desde: Italy	€45,00 +€19,35 envío €64,35 total	Añadir al carrito Detalles
	Les Beatles* - You've Got To Hide Your Love Away / Yesterday (7", Single)	Sello: Odeon	Cat. n.º: SO 10132	Estado del soporte: Near Mint (NM or M-) ☹	Condición de la funda: Very Good Plus (VG+)	COVER IS EXC Ver la página de la edición	owenscarlet ★★★★★ 99.0%, 1,913 valoraciones Enviado desde: France	€150,00 +€16,00 envío €166,00 total	Añadir al carrito Detalles
	Level 42 - Something About You (Sisa Mix) (12", EP)	Sello: Polydor	Cat. n.º: 883 362-1	Estado del soporte: Very Good (VG) ☹	Condición de la funda: Very Good (VG)	Ver la página de la edición	BrentSpar ★★★★★ 95.1%, 168 valoraciones Enviado desde: Netherlands	€2,25 +€12,00 envío €14,25 total	Añadir al carrito Detalles

Discos de vinilo a la venta recién añadidos en Discogs

De esta forma, un título adecuado sería. **“Singles y EPs de Michael Jackson”**

3. Descripción del dataset

En este dataset se obtendrá la lista completa de singles de Michael Jackson. En total, contendrá 294 registros, dónde se indicará información para cada single. Para cada single, extraemos su título, artista, año, país y sello discográfico. Guardaremos los datos en un archivo csv.

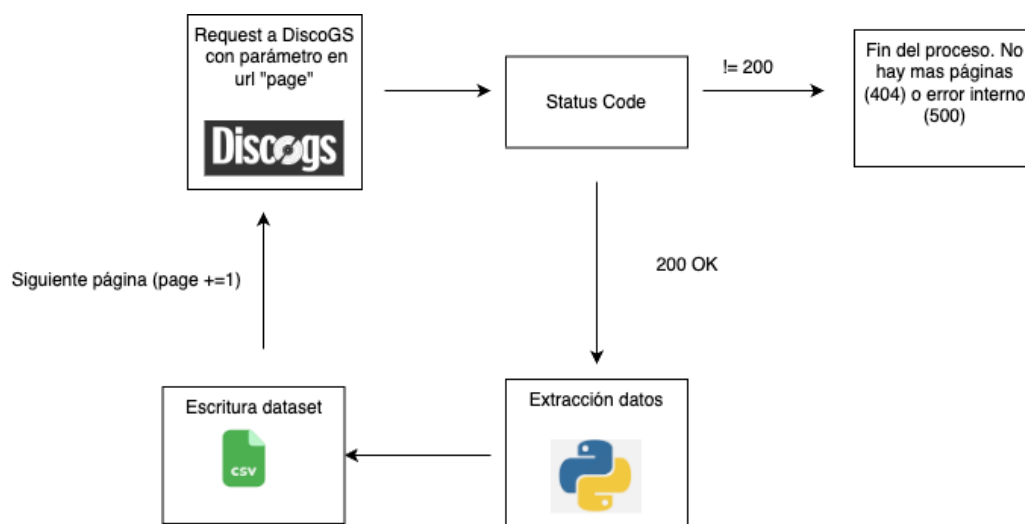
The screenshot shows the Discogs profile for Michael Jackson. It includes his name, birth and death dates, and a list of his releases. The discography section is highlighted, showing a list of singles and EPs with their respective labels and release years.

Releases	Count
Releases	607
Albums	21
Singles & EPs	294
Compilations	203

Página de dónde se obtendrán los datos

4. Representación gráfica

Esquema del funcionamiento del scrapper



5. Contenido

El dataset se guarda en formato CSV y contiene los siguientes campos, que indican información sobre cada uno de los singles.

Campo	Descripción	Tipo de datos	Ejemplo
title	Título del single	Texto	"Rock With You"
artist	Artista(s)	Texto	"Michael Jackson / Anita Ward"
year	Año de publicación	Número	1972
country	País de publicación	Texto	"Brazil"
label	Sello Discográfico	Texto	"Motown"

Contiene datos comprendidos entre 1972 y 2020. Aunque Michael Jackson falleció en 2009, se realizaron algunos lanzamientos de singles y colaboraciones de forma póstuma.

6. Propietario

No se han encontrado datasets similares a este. Además, al tratarse del listado de canciones de un artista, considero que los datos son públicos y no tienen propietario. Tendrán propietario los temas en sí, pero no un listado de todos ellos, que es lo que se obtiene en esta práctica.

Desde un punto de vista ético, no hallamos ningún inconveniente. Los datos extraídos únicamente se pueden emplear de forma lúdica y divulgativa, sin ningún fin comercial o militar.

7. Inspiración

Recientemente he recuperado una cadena de música con tocadiscos del trastero de mis padres, junto con una colección de vinilos de los años 80



Cadena de música como la mía

Esto ha despertado en mí un interés por el coleccionismo de música, concretamente de vinilos, con lo cuál suelo consultar bastante la web de **Discogs** para comprar vinilos.

Como aficionado a la música, me parece interesante poder obtener programáticamente una lista de creaciones de un artista. De la misma forma, se podría configurar el scrapper para obtener álbumes, compilatorios, grabaciones, etc.

Primeramente probé a scrapear vinilos de diversos artistas, pero obtenía un dataset demasiado pequeño (~10 registros en el mejor de los casos) y no me permitía implementar la navegación en el scrapper, ya que todo se mostraba en la misma página de resultados.

He escogido obtener Singles y EP de Michael Jackson porque se trata de un artista prolífico, ya que así se obtiene un buen dataset y puedo navegar entre páginas con el scraper

1 – 25 de 294 < Prev. 1 2 3 ... 11 12 Siguiente >

Páginas de resultados para la búsqueda de Singles y EPs de Michael Jackson

8. Licencia

Una licencia adecuada para estos datos sería **CCO: Public Domain License**.

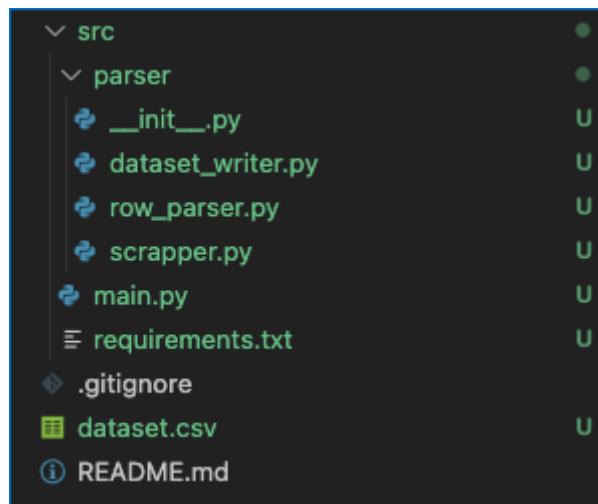
Esta licencia desiste de cualquier derecho de propiedad intelectual, liberando el documento al dominio público. Permite la copia, modificación, distribución y comunicación pública, incluso para fines comerciales, sin autorización previa.

Por otra parte, no se ofrece ninguna garantía respecto a la obra, y el autor renuncia a cualquier responsabilidad por el uso de la obra.

Considero que mi dataset “Singles y EPs de Michael Jackson” se ajusta a esta licencia, ya que los datos se encuentran fácilmente disponibles en internet y són de dominio público.

9. Código

El código se organiza mediante la siguiente estructura de archivos y carpetas



Clase Scraper

Realiza la petición para obtener el HTML de la página y convertirlo con BeautifulSoup.

Luego, itera a través de las páginas de resultados mediante la modificación del query param “page” de la URL.

Para cada página de resultados, escribe los datos en un CSV mediante la clase Dataset, que se muestra en el siguiente apartado.

Cuando el código de respuesta es diferente a 200, normalmente 404 para indicar que ya no hay más páginas, el scrapper se para de ejecutar.

Así mismo, muestra logs por pantalla sobre el progreso de scrapping.

El scrapper incluye las siguientes medidas para evitar ser bloqueado:

- Espera de 2 segundos entre peticiones

Se realiza entre llamadas a la página web mediante el comando `sleep(2)`

```
while True:
    page = self._fetch_page(page_num)

    if page.status_code != 200:
        print(f'Finished! Results: {self.total_results}')
        break

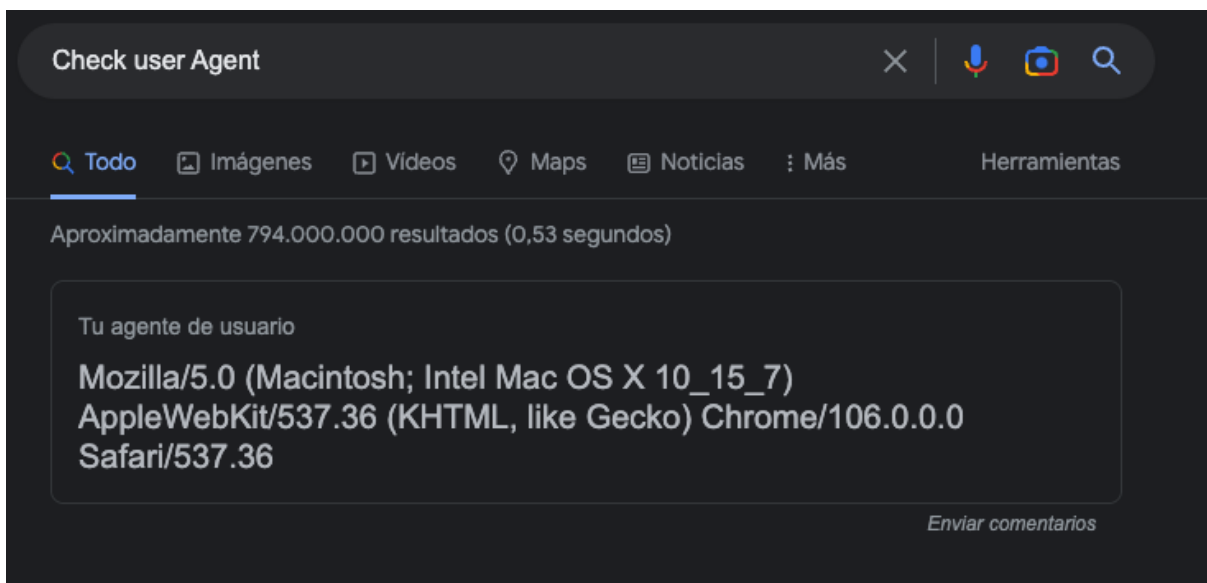
    print(f'--- Scrapping page {page_num} ---')

    self._scrape_page(page)
    page_num += 1
    sleep(2)

return
```

- Modificación del user-agent

Obtengo mi user-agent real y lo inserto en la llamada al sitio web



Clase DatasetWriter

Se encarga de crear el archivo del dataset y permite escribir filas. Utiliza la librería **csv** de Python.

```
src > parser > dataset_writer.py > ...
1  import csv
2
3
4  class DatasetWriter:
5      headers = ['title', 'artist', 'year', 'country', 'label']
6
7      def init(self):
8          with open('dataset.csv', 'w') as dataset:
9              writer = csv.writer(dataset)
10             writer.writerow(self.headers)
11
12     def write_row(self, row):
13         with open('dataset.csv', 'a') as dataset:
14             writer = csv.writer(dataset)
15             writer.writerow(row)
16
17
```

Clase TableRowParser

Clase que permite extraer la información de cada resultado de búsqueda mediante BeautifulSoup.

```
src > parser > row_parser.py > ...
1  class TableRowParser:
2      def __init__(self, tr):
3          self.tr = tr
4
5      def get_artist(self):
6          return self.tr.find('td', class_='artist').get_text()
7
8      def get_title(self):
9          title_data = self.tr.find('td', class_='title')
10         return title_data.find('a', recursive=False).get_text()
11
12     def get_label(self):
13         return self.tr.find('td', class_='label').get_text()
14
15     def get_country(self):
16         return self.tr.find('td', class_='country').get_text()
17
18     def get_year(self):
19         return self.tr.find('td', class_='year').get_text()
20
```

Uso de TableRowParser en Scrapper

```
for tr in results:
    parser = TableRowParser(tr)

    title = parser.get_title()
    artist = parser.get_artist()
    year = parser.get_year()
    country = parser.get_country()
    label = parser.get_label()

    self.dataset.write_row([title, artist, year, country, label])
```

Archivo main.py

Archivo principal, dónde se crea una nueva instancia de la clase Scrapper y se ejecuta el método **scrape()** para iniciar el proceso.

```
src > main.py > ...

1  from parser.scrapper import PageScrapper
2
3  scrapper = PageScrapper()
4
5  scrapper.scrape()
6  |
```

Cuando iniciamos el scrapper, se muestran los siguientes logs de progreso

```
➤ → scrapper git:(master) x /usr/bin/python3 /Users/daniel/Personal/scrapper/src/main.py
--- Scrapping page 1 ---
Results: 25
--- Scrapping page 2 ---
Results: 25
--- Scrapping page 3 ---
Results: 25
--- Scrapping page 4 ---
Results: 25
--- Scrapping page 5 ---
Results: 25
--- Scrapping page 6 ---
Results: 25
--- Scrapping page 7 ---
Results: 25
--- Scrapping page 8 ---
Results: 25
--- Scrapping page 9 ---
Results: 25
--- Scrapping page 10 ---
Results: 25
--- Scrapping page 11 ---
Results: 25
--- Scrapping page 12 ---
Results: 19
Finished! Results: 294
```


Así mismo, se ha generado el archivo dataset.csv

```
dataset.csv
1 title,artist,year,country,label
2 Ain't No Sunshine,Michael Jackson,1971,Jamaica,Tamla Motown
3 Maria (You Were The Only One),Michael Jackson,1971,Jamaica,"Tamla, Motown"
4 Got To Be There / Maria (You Were The Only One),Michael Jackson,1971,UK,Motown
5 We've Got A Good Thing Going = Una Buena Cosa Andando,Michael Jackson,1972,Venezuela,Motown
6 Jackson 5 Maxi,Michael Jackson / The Jackson 5,1972,Netherlands,"Tamla Motown, Tamla Motown"
7 You've Got A Friend / Ain't No Sunshine,Michael Jackson,1972,Venezuela,Tamla Motown
8 Ain't No Sunshine,Michael Jackson,1972,Netherlands,"Tamla Motown, Tamla Motown"
9 Ain't No Sunshine,Michael Jackson,1972,UK,Tamla Motown
10 Ben,Michael Jackson,1972,Brazil,Motown
11 Ain't No Sunshine / Ben,Michael Jackson,1972,Spain,"Tamla Motown, Tamla Motown"
12 I Wanna Be Where You Are.,Michael Jackson,1972,New Zealand,"Motown, Motown"
13 Rockin' Robin,Michael Jackson,1972,Australia,Motown
14 Ben,Michael Jackson,1972,Brazil,Motown
15 Rockin' Robin / Hey Big Brother,Michael Jackson / Rare Earth,1972,Italy,Tamla Motown
16 Got To Be There,Michael Jackson,1972,Brazil,Motown
17 Rockin' Robin,Michael Jackson,1972,Mexico,Tamla Motown
```

10. Dataset

El dataset ha sido subido a Zenodo y es accesible mediante la siguiente URL

- <https://zenodo.org/record/7225553>

Daniel Solá. (2022). Michael Jackson Singles & EPs [Data set]. Zenodo.
<https://doi.org/10.5281/zenodo.7225553>

11. Vídeo

El vídeo de la práctica se encuentra disponible en el siguiente enlace:

https://drive.google.com/file/d/1uOM0gk0DVFVNfhHp9EhZ2IGGqJxKsBzK/view?usp=share_link