

Extracción de datos con Python

Ejemplo 1:

El objetivo es obtener información de datos sobre la población del mundo utilizando la técnica de scraping o raspado. Los datos están disponibles en el sitio de Worldometer. Es un sitio de código abierto dirigido por un equipo internacional de desarrolladores e investigadores voluntarios, cuyo objetivo es poner las estadísticas globales a disposición de un amplio público en todo el mundo.

- Ejecutar cada una de las celdas de ejemplo.
- observar y revisar cada resultado.

```
In [ ]: 1 import csv
        2 import requests
        3 from bs4 import BeautifulSoup
        4 import pandas as pd
```

Crear una variable URL en formato string (texto) que contendrá el enlace a la página.

```
In [ ]: 1 url = "https://www.worldometers.info/world-population/population-by-country/"
```

Obtener los datos, utilizando la función `requests.get()`

```
In [ ]: 1 req = requests.get(url)
```

Con la clase `BeautifulSoup()`, extraer el código HTML de la página. En el argumento de esta función se selecciona el objeto `.text`

```
In [ ]: 1 soup = BeautifulSoup(req.text)
```

En los datos almacenados en la variable, buscar el código HTML con la palabra clave `"table"` y usar la función `.find_all()`

```
In [ ]: 1 data = soup.find_all("table")[0]
```

Leer el código HTML con el comando `.read_html(str())`, luego recuperar el primer y único elemento del objeto (el array)

```
In [ ]: 1 df_population = pd.read_html(str(data))[0]
```

Emitir los primeros elementos, utilizando la función `.head()` del `DataFrame`

```
In [ ]: 1 df_population.head()
```

Exportar los datos a un `.csv`

```
In [ ]: 1 export_csv = df_population.to_csv(r"archivo.csv", index = None, header = True)
```

Ejemplo 2:

El objetivo es hacer un raspado o scraping de citas de una página web y guardarlos en un archivo. Luego se leen y se emiten traducidas.

- Ejecutar el código a continuación.
- Revisar que se haya creado el archivo.
- Inspeccionar el archivo creado.
- Ejecutar el código siguiente.
- Inspeccionar el resultado.
- Quitar del texto las repeticiones de la comilla doble usando expresiones regulares.
- Investigar como guardar la última información obtenida en un archivo.

In []: ▶

```
1 # ejecutar
2 import requests
3 from bs4 import BeautifulSoup
4 import csv
5
6 def scrape_page(soup, quotes):
7     # recuperando todo el elemento HTML <div> en la página
8     quote_elements = soup.find_all('div', class_='quote')
9
10    # iterando sobre la lista de elementos de cotización
11    # para extraer los datos de interés y almacenarlos
12    # entre comillas
13    for quote_element in quote_elements:
14        # extrayendo el texto de la cita
15        text = quote_element.find('span', class_='text').text
16        # extrayendo el autor de la cita
17        author = quote_element.find('small', class_='author').text
18
19        # extrayendo la etiqueta <a> elementos HTML relacionados con la cita
20        tag_elements = quote_element.find('div', class_='tags').find_all('a', class_='tag')
21
22        # almacenar la lista de cadenas de etiquetas en una lista
23        tags = []
24        for tag_element in tag_elements:
25            tags.append(tag_element.text)
26
27        # agregando un diccionario que contiene los datos
28        # en un nuevo formato en la lista
29        quotes.append(
30            {
31                'text': text,
32                'author': author,
33                'tags': ', '.join(tags) # fusionar las etiquetas en una cadena "A, B, ..., Z"
34            }
35        )
```

```

In [ ]: ▶ 1 # La URL de La página de inicio del sitio web de destino
2 base_url = 'https://quotes.toscrape.com'
3
4 # definir el encabezado User-Agent para usar en la solicitud GET a continuación
5 headers = {
6 'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 \
7 (KHTML, like Gecko) Chrome/107.0.0.0 Safari/537.36'
8 }
9
10 # recuperar La página web de destino
11 page = requests.get(base_url, headers=headers)
12
13 # analizando La página web de destino con BeautifulSoup
14 soup = BeautifulSoup(page.text, 'html.parser')
15
16 # inicializando La variable que contendrá
17 # la lista de todos los datos
18 quotes = []
19
20 # raspando La página de inicio
21 scrape_page(soup, quotes)
22
23 # obteniendo el elemento "Siguiente →" HTML
24 next_li_element = soup.find('li', class_='next')
25
26 # si hay una página siguiente para raspar
27 while next_li_element is not None:
28     next_page_relative_url = next_li_element.find('a', href=True)['href']
29
30     # obteniendo La nueva página
31     page = requests.get(base_url + next_page_relative_url, headers=headers)
32
33     # analizando La nueva página
34     soup = BeautifulSoup(page.text, 'html.parser')
35
36     # raspando La nueva página
37     scrape_page(soup, quotes)
38
39     # buscando el elemento "Siguiente →" HTML en La nueva página
40     next_li_element = soup.find('li', class_='next')
41
42 # leer el archivo "citas.csv" y crearlo
43 # si no está presente
44 csv_file = open('citas.csv', 'w', encoding='utf-8', newline='')
45
46 # inicializando el objeto escritor para insertar datos
47 # en el archivo CSV
48 writer = csv.writer(csv_file)
49
50 # escribiendo el encabezado del archivo CSV
51 writer.writerow(['Text', 'Author', 'Tags'])
52
53 # escribiendo cada fila del CSV
54 for quote in quotes:
55     writer.writerow(quote.values())
56
57 # finalizar La operación y liberar los recursos
58 csv_file.close()

```

TextBlob es una biblioteca de Python para procesar datos textuales. Proporciona una API sencilla para el procesamiento del lenguaje natural (PLN).

```

In [ ]: ▶ 1 import textblob
2 from textblob import TextBlob

```

```

In [ ]: ▶ 1 file= open('citas.csv', encoding="utf8")
2 dataset= file.read()
3 trad=TextBlob(dataset)
4 traducido=trad.translate(from_lang='en',to='es')
5 traducido

```