



Índice

Índice.....	1
Introducción a Web Scraping.....	2
Pasos generales de Web Scraping.....	2
Tutorial de Web Scraping con Python, requests y BeautifulSoup.....	3
1. Hacer una Solicitud HTTP.....	3
2. Analizar el HTML.....	3
3. Seleccionar los Elementos.....	3
4. Extraer los Datos.....	3
5. Procesar y Almacenar los Datos.....	7
6. Iteración y Automatización.....	8



Introducción a Web Scraping

Web scraping es una técnica que permite extraer información de sitios web de forma automatizada. Puedes utilizar web scraping para obtener datos, como noticias, precios de productos, listas de empleos y mucho más. Sin embargo, es importante tener en cuenta que el web scraping debe hacerse de manera ética y respetar los términos de servicio de los sitios web a los que accedes.

Pasos generales de Web Scraping

Los pasos generales para realizar web scraping son los siguientes:

1. **Hacer una solicitud HTTP:** Utiliza la biblioteca requests para obtener el contenido HTML de la página web que deseas analizar.
2. **Analizar el HTML:** Utiliza la biblioteca BeautifulSoup para analizar el HTML y acceder a los datos de manera estructurada.
3. **Seleccionar los elementos:** Utiliza los métodos de BeautifulSoup, como find y find_all, para seleccionar los elementos que deseas extraer.
4. **Extraer los datos:** Accede al contenido de los elementos seleccionados y almacénalos para su posterior procesamiento.
5. **Procesar y almacenar los datos:** Procesa los datos extraídos y almacénalos en un formato útil, como CSV, JSON o una base de datos.
6. **Iteración y automatización:** Si es necesario, itera a través de múltiples páginas o sitios web para recopilar datos de manera automatizada.

Tutorial de Web Scraping con Python, requests y BeautifulSoup

1. Hacer una Solicitud HTTP

En este ejemplo, usaremos la página de noticias de BBC como sitio de prueba. Comencemos haciendo una solicitud HTTP para obtener el contenido HTML.

```
import requests

url = 'https://www.bbc.com/news'
response = requests.get(url)
```

2. Analizar el HTML

Usa BeautifulSoup para analizar el contenido HTML de la página:

```
from bs4 import BeautifulSoup

soup = BeautifulSoup(response.text, 'html.parser')
```

3. Seleccionar los Elementos

Vamos a seleccionar los títulos de las noticias en la página de la BBC. Inspecciona el HTML de la página para encontrar el selector adecuado. En este caso, los títulos están dentro de elementos `<h3>` con la clase "gs-c-promo-heading__title":

```
news_headlines = soup.find_all('h3', class_='gs-c-promo-heading__title')
```

4. Extraer los Datos

Itera a través de los elementos seleccionados y extrae el texto de los títulos de las noticias:

```
for headline in news_headlines:
    print(headline.text)
```



Aquí explicamos más detalles de los métodos `find` y `find_all` de `bs4`:

Método `find()`

El método `find()` se utiliza para encontrar la primera instancia de una etiqueta que cumple con ciertos criterios de búsqueda en un documento HTML. Puedes especificar qué etiqueta estás buscando y, opcionalmente, proporcionar atributos y contenido de texto para refinar tu búsqueda.

La sintaxis general es:

```
soup.find(name, attrs, text)
```

- **name**: El nombre de la etiqueta que deseas encontrar, como '`div`', '`p`', '`a`', etc.
- **attrs (opcional)**: Un diccionario de atributos y sus valores que deben coincidir para que se encuentre la etiqueta.
- **text (opcional)**: El contenido de texto que debe coincidir con el texto dentro de la etiqueta.

A continuación, se muestra un ejemplo de cómo usar `find()`:

```
from bs4 import BeautifulSoup

html = """
<html>
  <body>
    <div class="container">
      <p>Este es un párrafo</p>
      <p>Este es otro párrafo</p>
    </div>
  </body>
</html>
"""

soup = BeautifulSoup(html, 'html.parser')

# Encontrar la primera etiqueta 'p' dentro de la etiqueta 'div'
first_paragraph = soup.find('div').find('p')
print(first_paragraph.text)
```



```
<!DOCTYPE html>
<html>
<head>
    <title>Ejemplo con atributos</title>
</head>
<body>
    <div class="container">
        <p class="important">Este es un párrafo importante</p>
        <p>Este es un párrafo normal</p>
    </div>
</body>
</html>
```

Para encontrar el primer párrafo que tenga la clase **important**, puedes usar el método **find()** con el atributo **class**:

```
from bs4 import BeautifulSoup

html = """
<!DOCTYPE html>
<html>
<head>
    <title>Ejemplo con atributos</title>
</head>
<body>
    <div class="container">
        <p class="important">Este es un párrafo importante</p>
        <p>Este es un párrafo normal</p>
    </div>
</body>
</html>
"""

soup = BeautifulSoup(html, 'html.parser')

# Encontrar la primera etiqueta 'p' con la clase 'important'
important_paragraph = soup.find('p', class_='important')

print(important_paragraph.text)
```

En este ejemplo, `soup.find('p', class_='important')` busca la primera etiqueta p con la clase important. El resultado impreso será: “Este es un párrafo importante”.



Método find_all()

El método `find_all()` se utiliza para encontrar todas las instancias de una etiqueta que cumplan con ciertos criterios de búsqueda en un documento HTML. Puedes especificar la etiqueta que estás buscando y, opcionalmente, proporcionar atributos y contenido de texto para refinar tu búsqueda.

La sintaxis general es similar a la de `find()`:

```
soup.find_all(name, attrs, text)
```

- **name**: El nombre de la etiqueta que deseas encontrar, como '`div`', '`p`', '`a`', etc.
- **attrs (opcional)**: Un diccionario de atributos y sus valores que deben coincidir para que se encuentre la etiqueta.
- **text (opcional)**: El contenido de texto que debe coincidir con el texto dentro de la etiqueta.

A continuación, se muestra un ejemplo de cómo usar `find_all()`:

```
from bs4 import BeautifulSoup

html = """
<html>
  <body>
    <div class="container">
      <p>Este es un párrafo</p>
      <p>Este es otro párrafo</p>
    </div>
  </body>
</html>
"""

soup = BeautifulSoup(html, 'html.parser')

# Encontrar todas las etiquetas 'p' dentro de la etiqueta 'div'
paragraphs = soup.find('div').find_all('p')

for paragraph in paragraphs:
    print(paragraph.text)
```

En este ejemplo, `soup.find('div')` encuentra la etiqueta '`div`', y luego `find_all('p')` encuentra todas las etiquetas de párrafo dentro de esa etiqueta '`div`'. El resultado impreso es:

```
Este es un párrafo
Este es otro párrafo
```



5. Procesar y Almacenar los Datos

En este ejemplo, simplemente imprimimos los títulos en la consola. Sin embargo, en una aplicación real, podrías almacenar los datos en una lista, un archivo CSV o una base de datos.

Almacenar en diccionarios

Para almacenar la información recopilada durante el proceso de web scraping en una lista de diccionarios, puedes seguir estos pasos:

1. Define una lista vacía donde almacenarás los diccionarios.
2. Durante el proceso de web scraping, crea un diccionario para cada elemento de datos que deseas recopilar. Cada diccionario representará un conjunto de datos específico y contendrá pares clave-valor.
3. Agrega cada diccionario a la lista que definiste en el paso 1.

A continuación, te proporciono un ejemplo de cómo hacer esto. Supongamos que deseas recopilar información sobre libros desde un sitio web y almacenarla en una lista de diccionarios. Aquí está el código de ejemplo:

```
from bs4 import BeautifulSoup

# Ejemplo de HTML con información de libros
html = """
<html>
<head>
    <title>Libros</title>
</head>
<body>
    <div class="book">
        <h2>Título del Libro 1</h2>
        <p>Autor: Autor 1</p>
        <p>Género: Ficción</p>
    </div>
    <div class="book">
        <h2>Título del Libro 2</h2>
        <p>Autor: Autor 2</p>
        <p>Género: No Ficción</p>
    </div>
</body>
</html>
""""
```



```
soup = BeautifulSoup(html, 'html.parser')

# Define una lista para almacenar los diccionarios de libros
libros = []

# Encuentra todas las etiquetas 'div' con la clase 'book'
books = soup.find_all('div', class_='book')

# Itera a través de las etiquetas 'div' para recopilar información de
# cada libro
for book in books:
    titulo = book.find('h2').text
    autor = book.find('p', text='Autor:').find_next('p').text
    genero = book.find('p', text='Género:').find_next('p').text

    # Crea un diccionario para cada libro y agrega la información
    libro_info = {
        'Título': titulo,
        'Autor': autor,
        'Género': genero
    }

    # Agrega el diccionario a la lista de libros
    libros.append(libro_info)

# Imprime la lista de libros
for libro in libros:
    print(libro)
```

6. Iteración y Automatización

Si deseas recopilar datos de varias páginas o sitios web, puedes iterar a través de las URLs y repetir los pasos anteriores.