

Web scraping

Para economistas

Alejandro Acosta León

1 de junio de 2024

El material de este taller está disponible en:
<https://github.com/alejo-acosta/web-scraping>

Pueden descargar el material en la pestaña de código y luego en "Download ZIP".

O pueden clonar el repositorio con el siguiente comando:

```
git clone https://github.com/alejo-acosta/web-scraping
```

¿Qué esperas aprender en este taller?

Escríbanlo en el siguiente Link

Conceptos

- ¿Qué es web scraping?
- Problemas éticos, legales y mal uso del web scraping.
- Alternativas al web scraping.
- Herramientas y librerías necesarias.

Entorno

- Instalación de Python y librerías necesarias (requests, BeautifulSoup, pandas, selenium).
- Configuración del entorno de desarrollo (puede ser Jupyter Notebook, VSCode, etc.).

Nivel 1: data estructurada en HTML

- Pandas :)

Nivel 2: data semi-estructurada

- Inspección de elementos HTML.
- Extracción de datos con BeautifulSoup.
- Descarga de archivos.

Nivel 3: data no estructurada y que no usa HTML

- Automatización web con Selenium.
- Extracción de información relevante.

Imaginémonos que vivimos dentro del internet, en la gran nación democrática y soberana de Weblandia.

Así como en la vida real, Weblandia tiene bienes públicos y privados. Normalmente, el web scraping se hace sobre 'bienes públicos'.



N

VS

Los bienes públicos son no excluyentes y no rivales.
¿Pero realmente es así en Weblandia o incluso en la vida real?

¿Cuáles?

Ejemplos de uso del web scraping en la literatura

- Cavallo, Alberto, and Roberto Rigobon. 2016. "The Billion Prices Project: Using Online Prices for Measurement and Research." *Journal of Economic Perspectives*, 30 (2): 151-78.
- Glaeser, Edward L., Hyunjin Kim, and Michael Luca. 2018. "Nowcasting Gentrification: Using Yelp Data to Quantify Neighborhood Change." *AEA Papers and Proceedings*, 108: 77-82.

Manos a la obra, preparemos nuestro entorno de trabajo.