

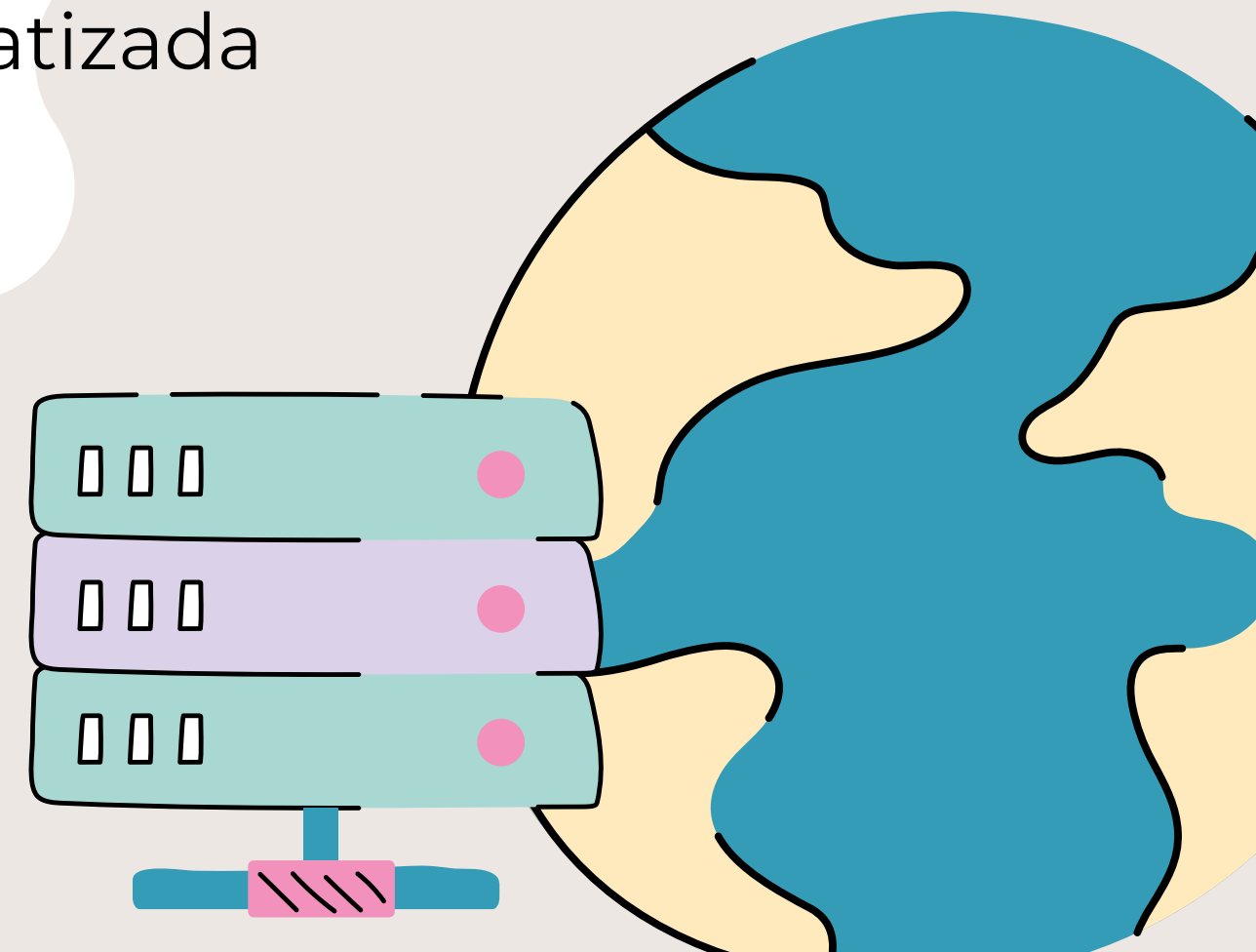


# Web Seraping con Python

Extracción de datos eficiente y automatizada



Por Juan Duran



# Introducción

El **Web Scraping** es el proceso de **recopilar información** de la **web** de manera programada, sin necesidad de introducir datos manualmente. Esta técnica es ampliamente utilizada en la industria para obtener **datos** de sitios web, realizar análisis de **tendencias**, **extraer** precios de productos, **recopilar** reseñas de usuarios y mucho más.

**Python** es uno de los lenguajes más populares para realizar Web Scraping debido a su versatilidad y la amplia variedad de librerías disponibles, como **BeautifulSoup** y **Scrapy**. Mientras que BeautifulSoup permite manipular el contenido HTML con facilidad, Scrapy es un framework más robusto que permite extracciones a gran escala. A pesar de su utilidad, es crucial considerar los aspectos **legales** y **éticos** del scraping. Muchos sitios web tienen políticas estrictas respecto a la extracción de datos, y el uso inadecuado de estas técnicas podría infringir sus términos de servicio.



# Puntos clave



## Automatización

- **Reduce** la necesidad de tareas manuales repetitivas.
- Permite **extraer** grandes volúmenes de información en poco tiempo.
- **Facilita** la recolección de datos estructurados para análisis.



## Herramientas

- **BeautifulSoup** permite una manipulación sencilla del código **HTML**.
- **Scrapy** facilita extracciones a gran escala mediante **spiders**.
- Ambas herramientas mejoran la **eficiencia** en la obtención de datos.



## Ética

- **Respetar** las políticas de acceso de cada sitio web.
- **Evitar** el scraping de información personal sin consentimiento.
- **Implementar** buenas prácticas para no sobrecargar servidores.

# Concepto de Web Scrapping

El **Web Scrapping** es una técnica utilizada para **extraer** datos de sitios web de forma **automática**. Consiste en enviar **solicitudes** a una página web, obtener el código fuente **HTML** y analizarlo para extraer la información deseada. A diferencia de la API de un sitio, que proporciona datos estructurados de manera oficial, el Web Scrapping permite obtener información de cualquier página, incluso cuando no hay una API disponible.

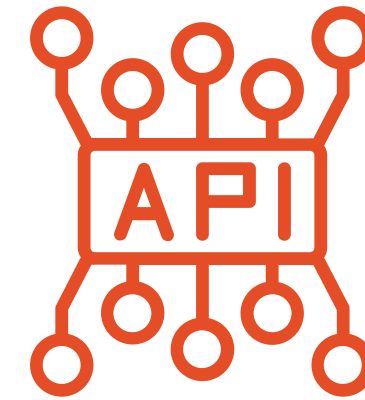
Las aplicaciones del Web Scrapping son diversas: desde el **monitoreo** de precios en e-commerce, **recopilación** de datos de investigación, hasta la **extracción** de noticias y **análisis** de redes sociales. Sin embargo, es fundamental realizar esta práctica de manera responsable y legal para evitar posibles penalizaciones.

**HTML**

**API**

**Automatización**

**Extracción**



**HTML**



# BeautifulSoup

**BeautifulSoup** es una librería de **Python** diseñada para facilitar la manipulación y **extracción** de datos de documentos **HTML** y **XML**. Proporciona herramientas intuitivas para navegar por el Document Object Model (**DOM**), lo que permite buscar etiquetas específicas, extraer información y modificar estructuras de datos con facilidad.

Para usar BeautifulSoup, primero se obtiene el código HTML de una página mediante la librería **requests** y luego se procesa con BeautifulSoup para encontrar los elementos deseados utilizando selectores **CSS** o funciones como `find()` y `find_all()`. Esto es especialmente útil para obtener tablas de datos, listas de artículos o información contenida en etiquetas específicas.

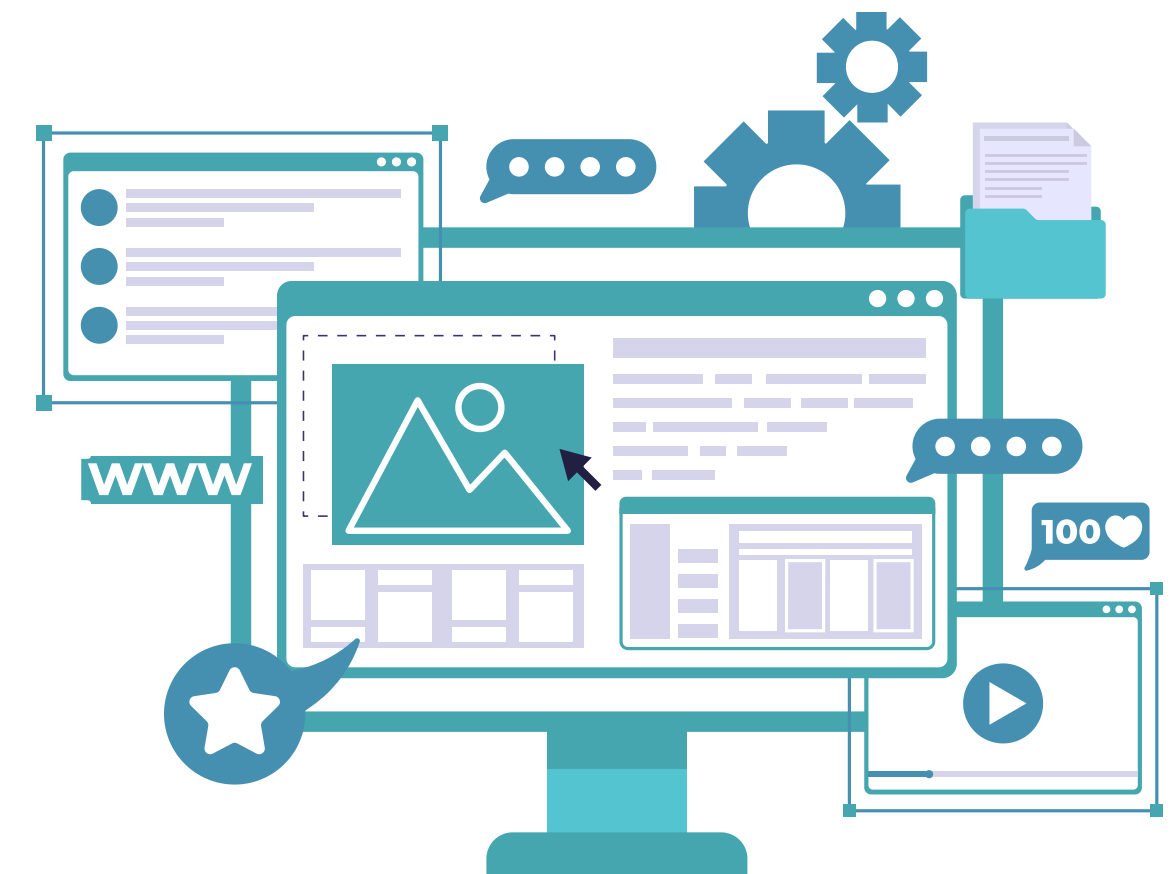
**HTML**

**Parsing**



**Filtrado**

**Navegación**



# Scrapy

**Scrapy** es un **framework** de Web Scraping más avanzado, diseñado para extraer datos de manera **eficiente** y **escalable**. Se basa en un sistema de **spiders**, que son programas que navegan automáticamente por la web siguiendo enlaces y extrayendo información según reglas predefinidas.

A diferencia de BeautifulSoup, Scrapy permite **gestionar** sesiones, **realizar** scraping de varias páginas de forma simultánea y **almacenar** datos de manera estructurada en bases de datos o archivos CSV. Además, incluye mecanismos para manejar **User-Agents** y evitar bloqueos por parte de los sitios web.

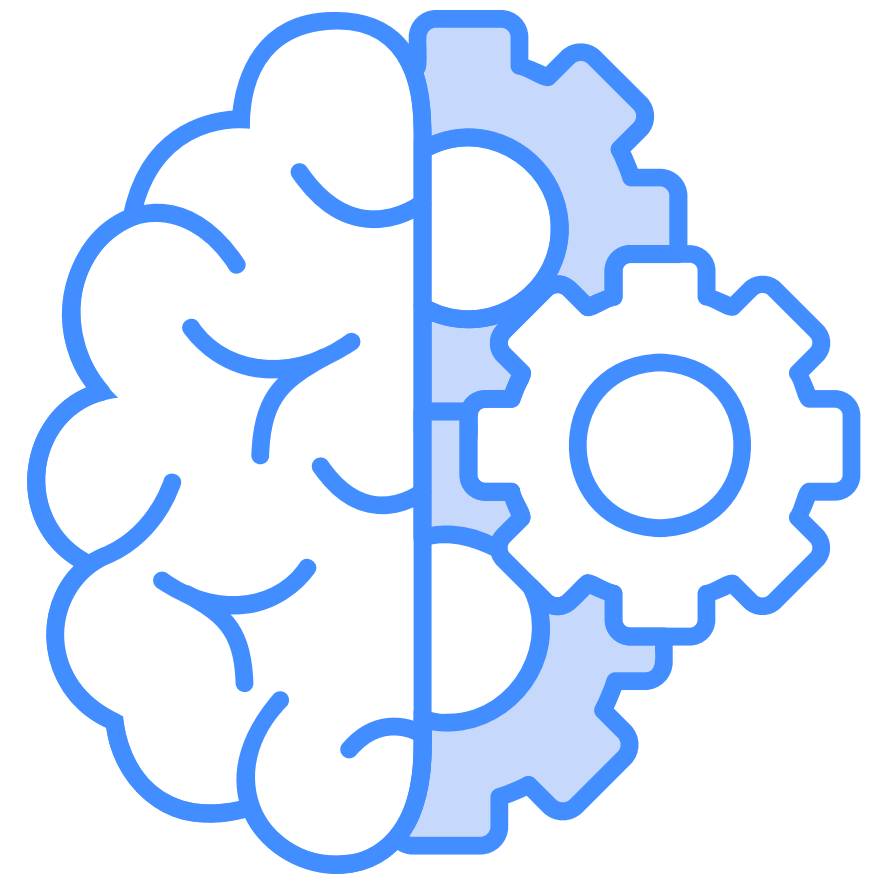
## Spiders



## Escalabilidad

## Framework

## Optimización



# Selenium

**Selenium** permite la **automatización** de navegadores para interactuar con sitios web que dependen de **JavaScript** y otros eventos dinámicos. Es especialmente útil cuando los datos a extraer no están disponibles en el HTML estático y requieren **simulaciones** de usuario, como desplazarse, hacer clic en botones o rellenar formularios.

Con Selenium, se pueden realizar pruebas de **navegación**, capturar datos después de la carga de scripts y **simular** acciones humanas en un entorno controlado. Es una herramienta clave para la automatización en **scraping** y **testing** de aplicaciones web.

**Automatización**

**Navegador**



**Interacción**

**JavaScript**

**{.js}**

**JavaScript**





# Conclusiones



## Web Scraping esencial

Una **herramienta** clave para recopilar información de la web de manera eficiente y automatizada

## BeautifulSoup y Scrappy

Son las principales **librerías** utilizadas para obtener y estructurar datos de sitios web.

## Selenium

Ideal para manipular sitios web **dinámicos** y simular interacciones de usuario.

## Consideraciones éticas

Siempre es necesario revisar las políticas de acceso de cada sitio para evitar problemas **legales**.

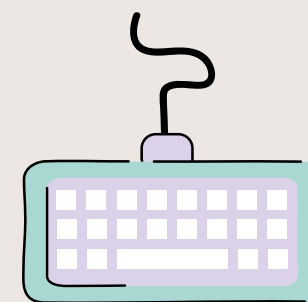
## Automatización

Permite ahorrar tiempo y mejorar la **eficiencia** en la extracción de datos.

## Habilidad valiosa

Conocer Web Scraping abre muchas **oportunidades** en el análisis y ciencia de datos.





# Gracias



Por Juan Duran

**“Coding, Gaming and Leveling Up”**