

Web scraping

Para economistas

Alejandro Acosta León

1 de junio de 2024

¿Qué esperas aprender en este taller?

Escríbanlo en el siguiente Jamboard

Conceptos

- ¿Qué es web scraping?
- Introducción al web scraping y sus aplicaciones.
- 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 que no usa HTML

- HTML y CSS.
- Inspección de elementos.
- Extracción de datos con BeautifulSoup.

¿Qué es Power BI?

- Power BI suele definirse como: Herramienta corporativa de autoservicio para hacer business intelligence.
- BI son herramientas que nos permiten extraer, conectar y visualizar datos.



<https://www.gartner.com/doc/reprints?id=1-24ZXJ0MU&ct=210107&st=sb>

¿Cuándo no usar Power BI?

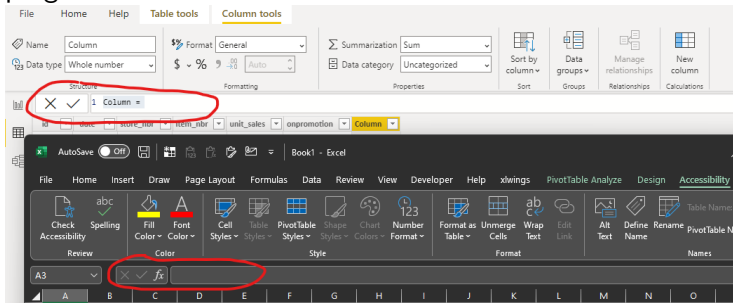
Power BI **no** es una herramienta de análisis o modelamiento de datos. Python, R, Stata, julia, etc. tienen una ventaja enorme en este campo. **El secreto es usar la herramienta adecuada para nuestro flujo de trabajo.**

Por ejemplo:

Extracción y manipulación datos	Análisis Exploratorio	Gráficos y reportes	Análisis Estadístico	Modelos
1.Python 2.Excel 3.Power BI 4.Stata	1.Python 2.Power BI 3.Excel 4.Stata	1.Python 2.Power BI 3.Excel	1.Python 2.Stata	1.Stata 2.R 3.Python

Los lenguajes de Power BI

- Power BI (al igual que Excel) tiene su propio "lenguaje de programación" llamado DAX.



- Sin embargo, también se utiliza otro lenguaje para extraer, transformar y cargar datos (ETL): M que es usado en Power Query.

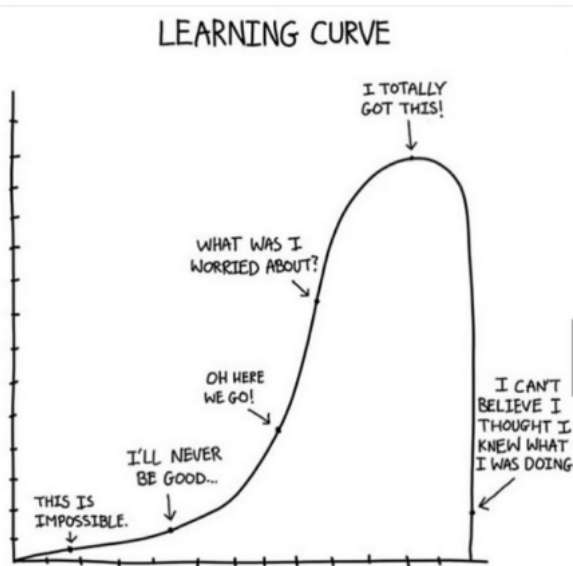
- Utilizaremos las ventas de Corporación la Favorita.
- Solo usaremos 1 % del total de la base.

<https://www.kaggle.com/c/favorita-grocery-sales-forecasting>

En esta sesión revisaremos:

- Refuerzo de tablas 'dimensiones' y modelo relacional de datos.
- Datos long vs wide.
- Métricas.
- Merge y append.
 - Funciones pivot y unpivot.
- Series de tiempo descriptivas (no predictivas).

Una pequeña motivación



Learning Curves - Popular Data Analysis Tools

