

# Proyecto Final: Automatización de Ingesta y Procesamiento de Datos Web hacia Base de Datos

## **Objetivo**

Desarrollar un flujo completo de ingestión de datos desde una fuente web hasta una base de datos relacional, aplicando técnicas avanzadas de web scraping, limpieza de datos, reglas de negocio y carga estructurada con SQLAlchemy. El resultado final debe representar un proceso automatizado, escalable y profesional, alineado con las buenas prácticas del curso.

## **Contexto del Proyecto**

La empresa requiere obtener información estructurada de libros del género Fantasy disponibles en el portal público Books to Scrape (<https://books.toscrape.com>). El objetivo es extraer esta información, procesarla según reglas definidas y almacenarla en una base de datos corporativa.

## **Requerimiento Principal**

Extraer todos los libros del género “Fantasy” y cargarlos en la base de datos cumpliendo las reglas de negocio.

## **Reglas del negocio**

La tabla final se llamará: books\_for\_sale

Debe contener al menos los siguientes campos:

- **book\_code:** Código único e incremental.
- **book\_name:** Nombre del libro (sin el contenido entre paréntesis).
- **book\_detail:** Texto extraído de los paréntesis del nombre (si aplica).
- **book\_price:** Precio del libro. Tipo numérico

## **Flujo del Proceso (Puntaje entre paréntesis)**

1. Navegar a la categoría Fantasy (1p)
2. Capturar al menos 3 atributos de un libro (2p)
3. Extraer la información de todos los libros de la primera página (2p)
4. Extraer los datos de todas las páginas (2p)
5. Crear la tabla books\_for\_sale usando SQLAlchemy ORM (2p)
6. Aplicar las reglas de negocio al cargar los datos (3p)
7. Implementar concurrencia en la carga hacia la base de datos (2p)

Adicionalmente se evaluará:

- No incluir librerías innecesarias (1p)
- Cerrar el navegador correctamente al terminar el scraping (1p)

Bono (Sólo complementará la participación):

- Desarrollo modular (No debe haber código suelto, todo en funciones) (2p)
- Aplicar Timing is everything (Uso de WebDriverWait) (2p)