

Criação de um Banco de Dados com SQLite após Web Scraping de Livros

Allyson de Almeida Sirvano

20 de maio de 2025

1 Introdução

Este trabalho tem como objetivo apresentar o processo completo de extração de dados do site *Books to Scrape* de forma automatizada utilizando a biblioteca Playwright. Os dados obtidos são armazenados em um banco de dados local SQLite, o que permite análises posteriores por meio de bibliotecas gráficas como Matplotlib e Seaborn.

2 Bibliotecas Utilizadas

2.1 Playwright

Responsável por controlar o navegador de forma programática, o Playwright permite navegar entre páginas, capturar elementos e simular interações humanas como cliques e rolagens.

Listing 1: Importação do Playwright

```
1 from playwright.sync_api import sync_playwright
```

O método `sync_playwright()` inicia uma instância sincronizada do navegador, permitindo o uso das funções do Playwright em modo não assíncrono, facilitando o uso em scripts simples.

2.2 SQLite3

Utilizado para a criação de um banco de dados leve, local e sem dependências externas. SQLite3 permite a persistência dos dados extraídos.

Listing 2: Importação do SQLite

```
1 import sqlite3
```

A biblioteca `sqlite3` provê uma interface para interagir com bancos de dados SQLite por meio de comandos SQL embutidos em Python.

2.3 Matplotlib e Seaborn

Estas bibliotecas são usadas para visualização de dados. O Matplotlib oferece maior controle sobre os gráficos, enquanto o Seaborn simplifica a criação de visualizações estatísticas.

Listing 3: Importação para visualização

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
```

3 Processo de Extração de Dados

A navegação é feita de forma automatizada, capturando os links de todos os livros presentes na página inicial. Posteriormente, cada link é acessado individualmente para coletar informações detalhadas.

Listing 4: Extração dos links dos livros

```
1 livros_links = page.eval_on_selector_all(
2     '.product_pod h3 a',
3     '(elements) => elements.map(element => ({ title: element.
4         getAttribute("title"), href: element.href })))'
```

A função `eval_on_selector_all` executa um código JavaScript dentro do contexto da página, coletando todos os elementos que correspondem ao seletor CSS informado. Neste caso, obtém-se os títulos e URLs dos livros exibidos na página.

4 Tratamento dos Dados

Os dados obtidos são transformados para os tipos corretos (números, texto) e validados. Registros com informações incorretas ou ausentes são descartados para garantir a integridade do banco.

Listing 5: Validação dos dados

```
1 if not titulo or preco <= 0 or quantidade < 0 or avaliacao not in
   range(1, 6):
2     print(f"Dados inválidos encontrados: {livro}")
3     continue
```

A condição verifica se o título é vazio, se o preço é menor ou igual a zero, se a quantidade é negativa ou se a avaliação não está entre 1 e 5 estrelas. Dados que falharem essas validações são ignorados.

5 Criação do Banco de Dados

Duas tabelas principais foram criadas: uma para categorias e outra para livros. Cada livro está vinculado à sua categoria por meio de uma chave estrangeira.

Listing 6: Criação das tabelas

```
1 CREATE TABLE IF NOT EXISTS categorias (  
2     id INTEGER PRIMARY KEY AUTOINCREMENT,  
3     nome TEXT NOT NULL UNIQUE  
4 );  
5 CREATE TABLE IF NOT EXISTS livros (  
6     id INTEGER PRIMARY KEY AUTOINCREMENT,  
7     titulo TEXT NOT NULL UNIQUE,  
8     preco REAL NOT NULL,  
9     quantidade INTEGER NOT NULL,  
10    avaliacao INTEGER NOT NULL,  
11    categoria_id INTEGER,  
12    FOREIGN KEY (categoria_id) REFERENCES categorias(id)  
13 );
```

CREATE TABLE IF NOT EXISTS cria as tabelas apenas se elas ainda não existirem. Campos como AUTOINCREMENT garantem que o ID seja gerado automaticamente.

6 Inserção dos Dados no Banco

A inserção é feita evitando duplicidades. Cada categoria é verificada previamente e, se ainda não existir, adicionada ao banco. Em seguida, os livros são inseridos com referência à categoria correspondente.

Listing 7: Inserção de livros com categoria associada

```
1 INSERT INTO livros (titulo, preco, quantidade, avaliacao,  
2     categoria_id)  
VALUES (?, ?, ?, ?, ?)
```

A sintaxe com ? permite uso de parâmetros seguros com o SQLite, evitando injeção de SQL.

7 Reorganização das Categorias

Para manter os IDs das categorias organizados, foi criada uma tabela temporária com nova sequência. Os registros da tabela de livros foram atualizados para refletir os novos IDs.

Listing 8: Atualização dos IDs das categorias

```
1 UPDATE livros  
2 SET categoria_id = (  
3     SELECT id FROM categorias_temp  
4     WHERE categorias_temp.nome = (  
5         SELECT nome FROM categorias  
6         WHERE categorias.id = livros.categoria_id  
7     )  
8 )
```

Este comando substitui os IDs antigos pelos novos baseando-se no nome da categoria. Utiliza subqueries aninhadas para buscar a equivalência entre as tabelas.

8 Indicadores de Performance

Foram geradas consultas para identificar padrões relevantes, como livros com alta avaliação, estoques baixos e preços médios.

Listing 9: Consultas para indicadores

```
1 SELECT COUNT(*) FROM livros WHERE avaliacao >= 4
2 SELECT COUNT(*) FROM livros WHERE quantidade <= 5
3 SELECT AVG(preco) FROM livros WHERE avaliacao >= 4
```

Cada comando SQL realiza uma análise específica: contagem de livros bem avaliados, livros com estoque baixo e média de preços dos livros de boa avaliação.

9 Visualização de Dados

Os dados foram visualizados por meio de gráficos de barras e pizza. O gráfico de barras exibe a quantidade de livros por nota de avaliação.

Listing 10: Gráfico de barras com Seaborn

```
1 sns.barplot(
2     x=distribuicao.keys(),
3     y=distribuicao.values(),
4     palette="muted",
5     ax=axes[0]
6 )
```

A função `barplot()` do Seaborn recebe as chaves (notas) e valores (quantidades), e renderiza um gráfico de barras com paleta de cores especificada.

O gráfico de pizza apresenta a proporção percentual de livros em cada faixa de avaliação.

Listing 11: Gráfico de pizza com Matplotlib

```
1 axes[1].pie(
2     percentual,
3     labels=[f'{p} estrelas ({v:.1f}%)' for p, v in zip(distribuicao
4     .keys(), percentual)],
5     startangle=140,
6     colors=sns.color_palette("pastel", len(distribuicao))
7 )
```

O método `pie()` do Matplotlib gera um gráfico circular usando os percentuais e rótulos personalizados. O parâmetro `startangle` gira o início do gráfico e as cores são definidas com a paleta "pastel" do Seaborn.