

Extract Structured Data of Books from Amazon and Barnes & Noble

Lan Bai, Chaoqun Mei, Yuzhe Ma

March 22, 2018

Abstract In this project stage, we used the open-source tool Scrapy to extract structured data of books from Amazon and Barnes & Noble. For each book, the publisher, time, author, and title attributes were extracted as a tuple. We finally got two tables in CSV format, one has 3079 tuples for Amazon books, and the other has 3179 tuples for Barnes & Noble books.

1. Data Sources

We selected Amazon books and Barnes & Noble books as the two Web data sources for this project. Since there are billions of books on Amazon books and Barnes & Noble books, we narrowed down to fiction, health, and business books on Amazon books, and history, health, business, and science books on Barnes & Noble books. For each book, the publisher, time, author, and title were extracted as a tuple.

2. Data Extraction

We extracted the structured data of books from Amazon and Barnes & Noble using Scrapy. Scrapy is a free and open source web crawling framework, written in Python. We located the information we would like to extract by searching the web source code, then the publisher, time, author, and title were extracted as a tuple for each book page by page. If there was a break during the extracting, we continued the extraction from a new start URL.

3. Type of Entity and Description of the Two Tables

The entity we extracted is book, we extracted the information of 3079 books from Amazon books, and 3179 from Barnes & Noble books. For each of the two Web data sources, we finally got a table in CSV format. There are 3079 tuples in the table for Amazon books, and 3179 tuples in the table for Barnes & Noble books. Each tuple has four attributes, namely, publisher, time, author, and title.

4. Tool

We used the open-source tool Scrapy in this project stage. Scrapy is an application framework written in Python for crawling web sites and extracting structured data which can be used for a wide range of useful applications, like data mining, information processing or historical archival. Even though Scrapy was originally designed for web scraping, it can also be used to

extract data using APIs (such as Amazon Associates Web Services) or as a general purpose web crawler.