

Final Project Data Checkpoint

----- Xinyi Zhao

Project code

https://github.com/sarahzhao21/final_project.git

Data sources

1. The New York Times Best Sellers: <https://www.nytimes.com/books/best-sellers/>

This website provides the lists of the best seller books in different categories in "New York Times".
Format: HTML, JSON

How: I did web crawling and scraping to get the information of best seller books and get the url of another link in "Apple Books" for each book and do the web crawling to "Apple Books".

Summary of data:

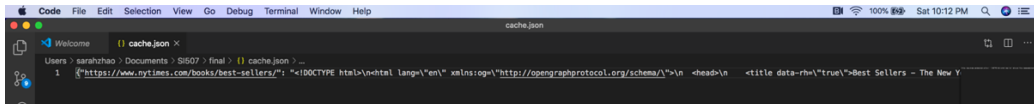
- 1) There are 140 records available (11 book categories and 10 to 15 best seller books for each categories)
- 2) 140 records have been retrieved.
- 3) description of records:

fields	Data type	Description
Id	INTEGER	The Id number of the book
Title	TEXT	The title of the book
Category	TEXT	What kind of book it is (fiction, novel)
Author	TEXT	The author who wrote the book
Publisher	TEXT	The publisher of the book
Rank	INTEGER	The rank of this book in it's category last week
Weeks_on_the_list	INTEGER	How many weeks the book was on the best-seller list
Description	TEXT	A short review about this book
Apple_url	TEXT	The url that connects to the website of 'Apple Books'

- 4) Evidence of caching

```
6 from bs4 import BeautifulSoup
7 import requests
8 import json
9 import sqlite3
10
11 CACHE_FILE_NAME = 'cache.json'
12 CACHE_DICT = {}
13
14 def load_cache():
15     """ Opens the cache file if it exists and loads the JSON into
16     the CACHE_DICT dictionary.
17     If the cache file doesn't exist, creates a new cache dictionary
18
19     Parameters
20     -----
21     None
22
23     Returns
24     -----
25     The opened cache: dict
26     """
27     try:
28         cache_file = open(CACHE_FILE_NAME, 'r')
29         cache_file_contents = cache_file.read()
30         cache = json.loads(cache_file_contents)
31         cache_file.close()
32     except:
33         cache = {}
34     return cache
35
36 def save_cache(cache):
37     """ Saves the current state of the cache to disk
38
39     Parameters
40     -----
41     cache_dict: dict
42         The dictionary to save
43
44     Returns
45     -----
46     None
47     """
48     cache_file = open(CACHE_FILE_NAME, 'w')
49     contents_to_write = json.dumps(cache)
50     cache_file.write(contents_to_write)
51     cache_file.close()
```

```
53
54 def make_url_request_using_cache(url, cache):
55     """ Check the cache for a saved result for this baseurl+params:values
56     combo. If the result is found, return it. Otherwise send a new
57     request, save it, then return it.
58
59     Parameters
60     -----
61     url: string
62         The basic URL of the website
63     cache: dict
64         A dictionary where the url_text saved
65
66     Returns
67     -----
68     url_text
69         the result of the query as a url_text get from the website
70     """
71     if url in cache.keys(): # the url is our unique key
72         print("Using cache")
73         return cache[url]
74     else:
75         print("Fetching")
76         response = requests.get(url)
77         cache[url] = response.text
78         save_cache(cache)
79         return cache[url]
80
81
82 def get_booksite_dic(base_url):
83     url_text = make_url_request_using_cache(base_url, CACHE_DICT)
84     soup = BeautifulSoup(url_text, 'html.parser')
85     booksite_dic = {}
86     bookgroup = soup.find('div', itemtype="http://schema.org/ItemList")
87     booksites = bookgroup.find_all('div', class_="css-v2K15d")
88     for booksite in booksites:
89         path = booksite.find('h2').find('a')['href']
90         site_url = 'https://www.nytimes.com' + path
91         key = booksite.find('h2')['id']
92         booksite_dic[key] = site_url
93     return booksite_dic
94
95 def get_bookinfo_list(booksite_dic):
96     book_list = []
97     i = 1
98     for key,value in booksite_dic.items():
99         url_text = make_url_request_using_cache(value, CACHE_DICT)
```



2. Apple Books

Each book from the New York Time Best Sellers has an 'Apple Books' url, here is one of the example: <https://books.apple.com/us/book/american-dirt-oprahs-book-club/id1459876569>

This website provides all kinds of detailed information about the books showed on New York Times.

Format: HTML, JSON

How: By accessing the beautiful soup of the Apple Books url for each book, I did the web scraping to obtain the information of 'title', 'genre', 'price', 'rating', 'released date', 'language', 'length', 'seller' and 'size'.

Summary of data:

- 1) There are 140 records available (each record refers to one book on New York Time best seller)
- 2) 140 records have been retrieved.
- 3) description of records:

fields	Data type	Description
Id	INTEGER	The Id number of the book
Title	TEXT	The title of the book
Rating	REAL	The rating of the book on "Apple Books"
Price	REAL	The price of the book on "Apple Books"
Genre	TEXT	The genre of the book on "Apple Books"
Released date	TEXT	The date that the book released on the website
Language	INTEGER	What kind of language the book was written by
Length	INTEGER	The page number of the book
Seller	TEXT	The seller of the book
Size	REAL	How many 'MB' of the book

4) Evidence of caching

```

Code File Edit Selection View Go Debug Terminal Window Help
Welcome final_proj.py
Users > sarahzhao > Documents > S1607 > final > final_proj.py > ...
1 6 from bs4 import BeautifulSoup
7 import requests
8 import json
9 import sqlite3
10
11 CACHE_FILE_NAME = 'cache.json'
12 CACHE_DICT = {}
13
14 def load_cache():
15     """ Opens the cache file if it exists and loads the JSON into
16     the CACHE_DICT dictionary.
17     If the cache file doesn't exist, creates a new cache dictionary
18
19     Parameters
20     -----
21     None
22
23     Returns
24     -----
25     The opened cache dict
26     """
27     try:
28         cache_file = open(CACHE_FILE_NAME, 'r')
29         cache_file_contents = cache_file.read()
30         cache = json.loads(cache_file_contents)
31         cache_file.close()
32     except:
33         cache = {}
34     return cache
35
36 def save_cache(cache):
37     """ Saves the current state of the cache to disk
38
39     Parameters
40     -----
41     cache_dict: dict
42         The dictionary to save
43
44     Returns
45     -----
46     None
47     """
48     cache_file = open(CACHE_FILE_NAME, 'w')
49     contents_to_write = json.dumps(cache)
50     cache_file.write(contents_to_write)
51     cache_file.close()

```

```

Code File Edit Selection View Go Debug Terminal Window Help
Welcome final_proj.py
Users > sarahzhao > Documents > S1607 > final > final_proj.py > ...
54 def make_url_request_using_cache(url, cache):
55     """ Check the cache for a saved result for this baseurl+params+values
56     combo. If the result is found, return it. Otherwise send a new
57     request, save it, then return it.
58
59     Parameters
60     -----
61     url: string
62         The basic URL of the website
63     cache: dict
64         A dictionary where the url_text saved
65
66     Returns
67     -----
68     url_text
69         the result of the query as a url_text get from the website
70     """
71     if url in cache.keys(): # the url is our unique key
72         print("using cache")
73         return cache[url]
74     else:
75         print("fetching")
76         response = requests.get(url)
77         cache[url] = response.text
78         save_cache(cache)
79         return cache[url]
80
81 def get_booksite_dic(base_url):
82     url_text = make_url_request_using_cache(base_url, CACHE_DICT)
83     soup = BeautifulSoup(url_text, 'html.parser')
84     bookgroup = soup.find('div', itype="http://schema.org/ItemList")
85     booksites = bookgroup.find_all('div', class_="css-v2k15P")
86     for booksite in booksites:
87         path = booksite.find('h2').find('a')['href']
88         site_url = "https://www.nytimes.com" + path
89         key = booksite.find('h2')['id']
90         booksite_dic[key] = site_url
91     return booksite_dic
92
93 def get_bookinfo_list(booksite_dic):
94     book_list = []
95     i = 1
96     for key, value in booksite_dic.items():
97         url_text = make_url_request_using_cache(value, CACHE_DICT)
98         soup = BeautifulSoup(url_text, 'html.parser')

```

```

Code File Edit Selection View Go Debug Terminal Window Help
Welcome final_proj.py
Users > sarahzhao > Documents > S1607 > final > final_proj.py > ...
171
172 def get_applebook_list(book_list):
173     apple_urls = []
174     apple_list = []
175     for book in book_list:
176         apple_urls.append(book[-1])
177
178     i = 1
179     for url in apple_urls:
180         apple_row = []
181         url_text = make_url_request_using_cache(url, CACHE_DICT)
182         soup = BeautifulSoup(url_text, 'html.parser')
183         try:
184             title = soup.find('h1', class_="product-header__title book-header_
185             title = None
186         except:
187             rating = soup.find('p', class_="we-rating-count star-rati
188             rating = None
189         except:
190             price = soup.find('li', class_="inline-list-item inline-list_ite
191             price = None
192         except:
193             section = soup.find('section', class_="l-content-width l-row l-row-pe
194             figure = section.find_all('figure')
195             figure = None
196         if figure is None:
197             genre = None
198             release_date = None
199             language = None
200             length = None
201             seller = None
202             size = None
203         else:
204             genre = figure[0].find('div', class_="book-badge__caption").text.s
205             release_date1 = figure[1].find('div', class_="book-badge__value").
206             release_date2 = figure[1].find('div', class_="book-badge__caption")
207             release_date = build_date(release_date1, release_date2)
208             language = figure[2].find('div', class_="book-badge__caption").tes
209             try:
210                 length = figure[3].find('div', class_="book-badge__value").tes
211             except:
212                 length = None
213

```

Database

A file named 'best_seller_book.sqlite' has been created by the program 'final_proj.py', when it was opened by DB Browser, there are two tables in this file: 'Best_seller' and 'Apple_book':

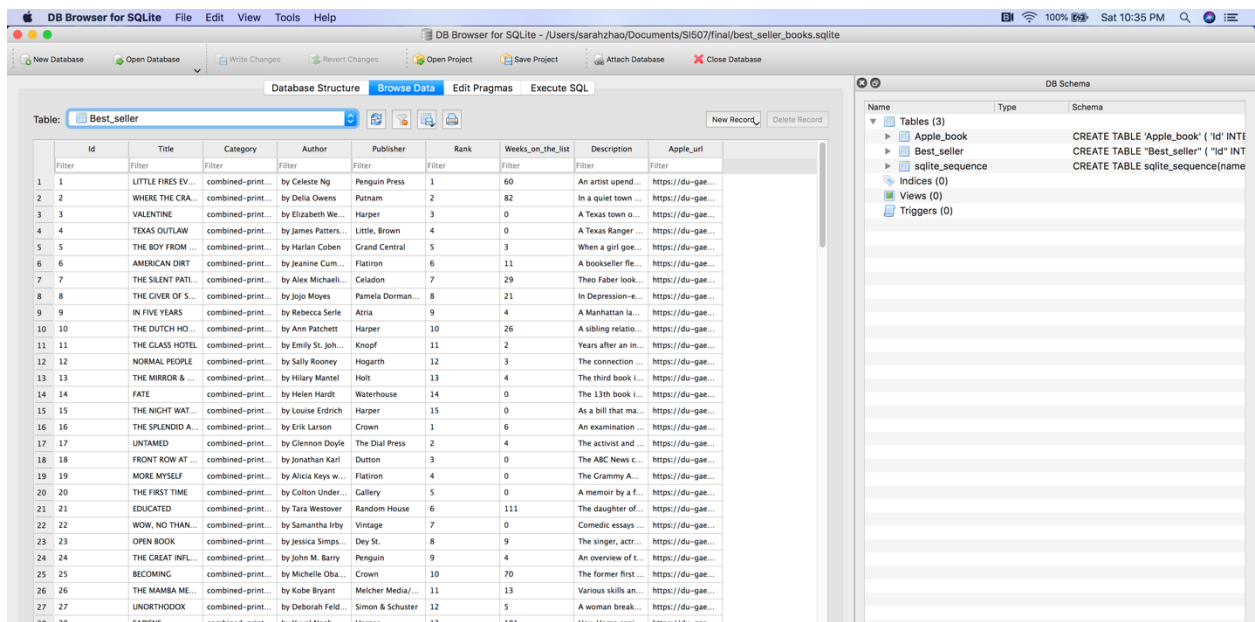
1. 'Best_seller' has 9 columns and 140 rows. Each row represents the information of a book, the information of the columns is listed here:

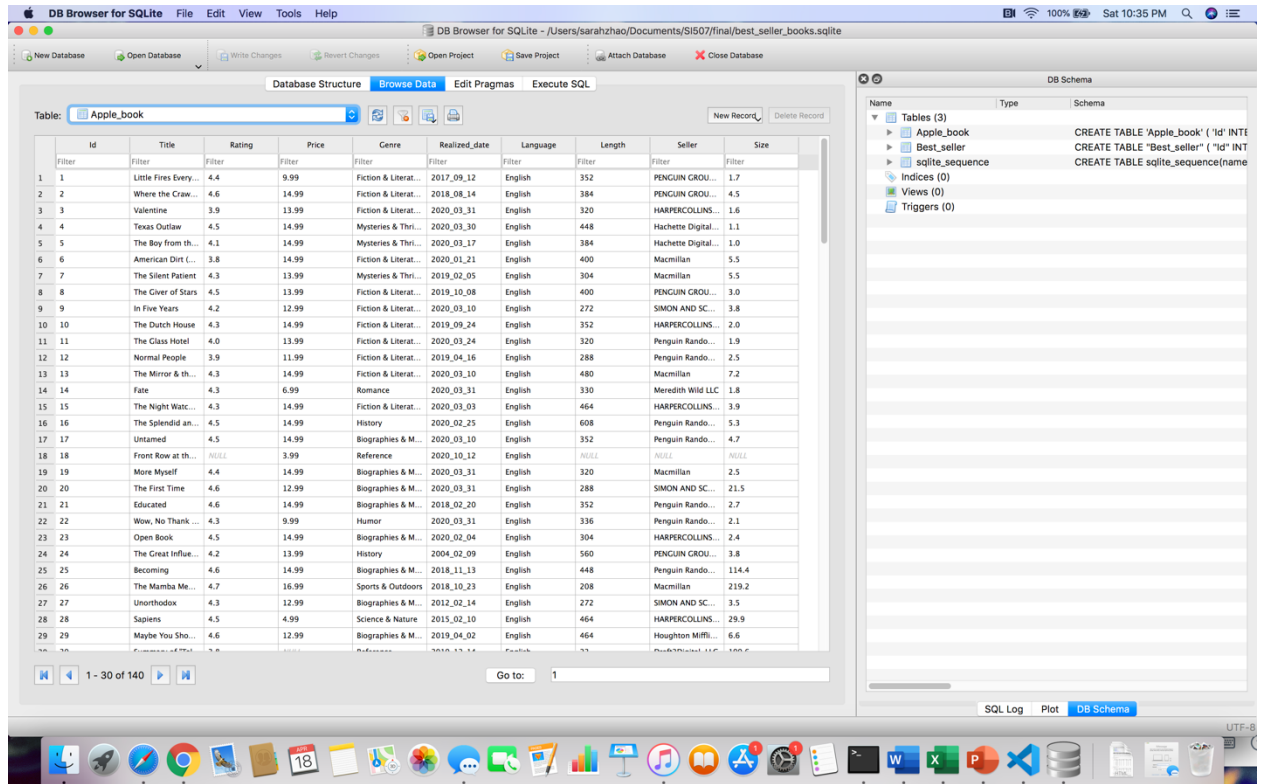
Column name	Data type
Id	INTEGER
Title	Text
Category	TEXT
Author	TEXT
Publisher	TEXT
Rank	INTEGER
Weeks on the list	INTEGER
Description	TEXT
Apple_url	TEXT

2. 'Apple_book' has 10 columns and 140 rows. Each row represents the information of a book on "Apple Books" and refer to the same book on the table 'Best_seller' with the same order. The information of the columns is listed here:

Column name	Data type
Id	INTEGER
Title	TEXT
Rating	REAL
Price	REAL
Genre	TEXT
Released date	TEXT
Language	INTEGER
Length	INTEGER
Seller	TEXT
Size	REAL

3. The primary key is the 'Id' on 'Best_seller' table and the foreign key is the 'Id' on 'Apple_book' table. The sequences of these two sets of 'Id's are the same.
4. Screen shot:





Interaction and Presentation Plans

Approach: Flask, Plotly

I plan to use Flask to design a submission form for the users to select what kinds of books they want to know. The form for users should look like this:

Welcome to the New York Time Best Seller Book Brower!

Select Category:

- Combined Print & E-Book Fiction
- Combined Print & E-Book Nonfiction
- Hardcover Fiction
- Hardcover Nonfiction
- Paperback Trade Fiction
- Paperback Trade Nonfiction
- Advice, How-To & Miscellaneous
- Children's Middle Grade Hardcover
- Children's Picture Books
- Children's Series
- Young Adult Hardcover

Sort by:

- Rank

- Rating on Apple Books
- Price on Apple Books
- Length on Apple Books

Sort direction:

- High to low
- Low to high

Plot results?

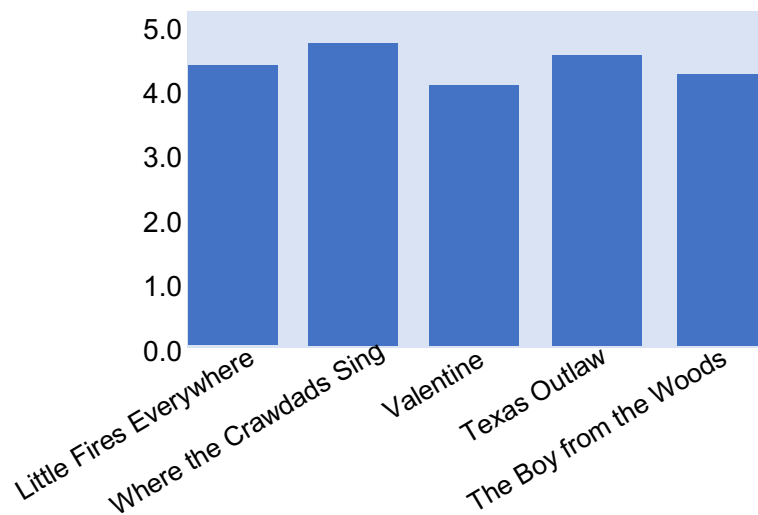
The presentation to users should look like this:

Here are your results:

Title	Genre	Author	Rank	Rating	Price	Length	Seller
Little Fires Everywhere	Fiction & Literature	Celeste Ng	1	4.4	9.99	352	PENGUIN GROUP
Where the Crawdads Sing	Fiction & Literature	Delia Owens	2	4.6	14.99	384	PENGUIN GROUP
Valentine	Fiction & Literature	Elizabeth Wetmore	3	3.9	13.99	320	HARPERCOLLINS
Texas Outlaw	Mysteries & Thrillers	James Patterson	4	4.5	14.99	448	Hachette Digital, Inc
The Boy from the Woods.	Mysteries & Thrillers	Harlan Coben	5	4.1	14.99	384	Hachette Digital, Inc

Or this:

Here is your graph:



The values that the graph will show depend on the 'Sort by' selected by the users. For example, if the user selects 'Sort by' 'Rating on Apple Books', the bar chart will show the ratings of the books labeled by the title of each book.