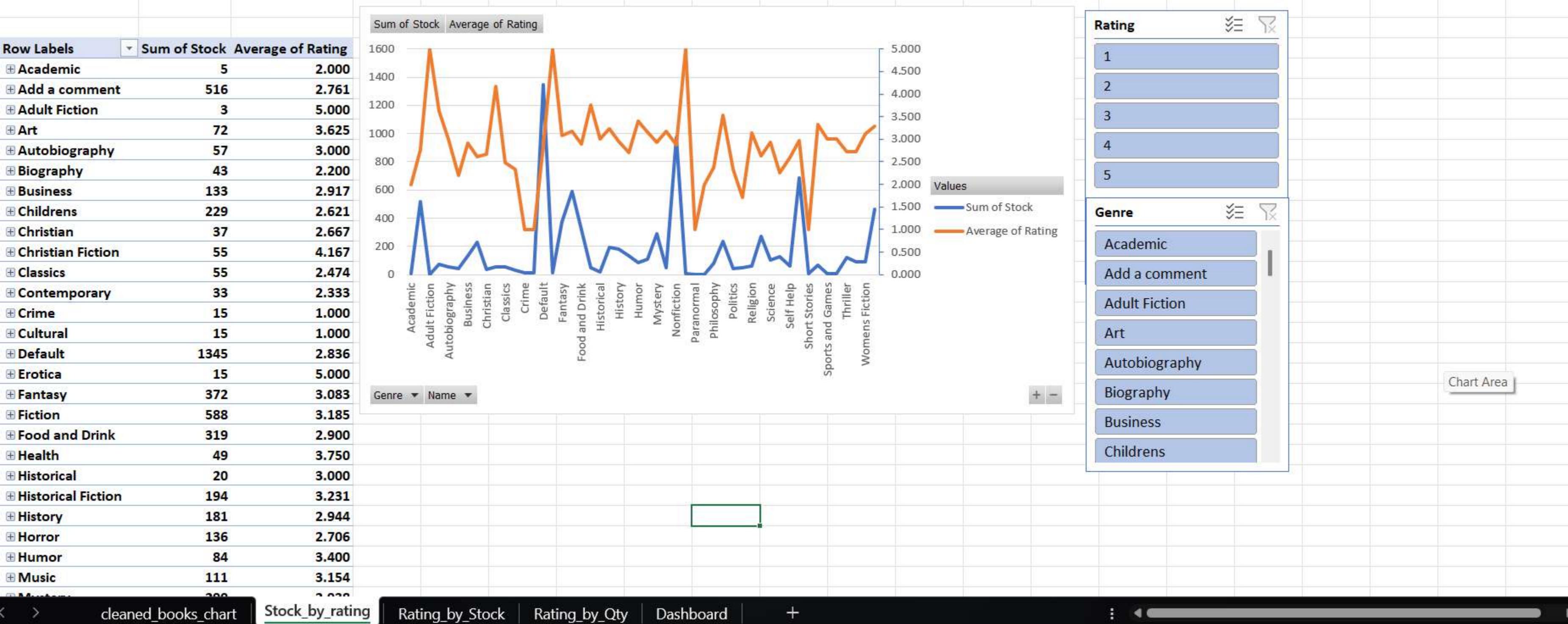


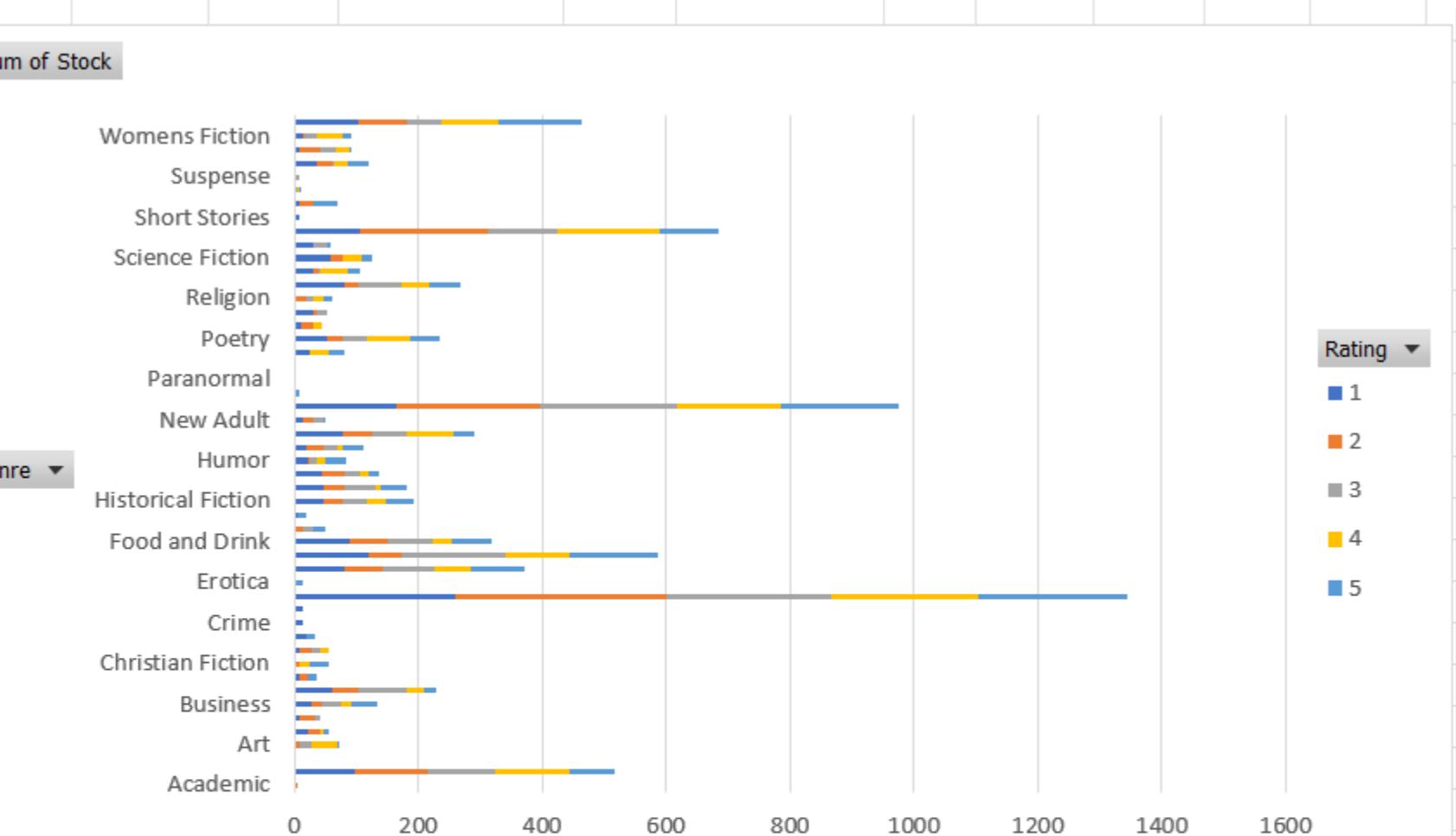
```
1 from bs4 import BeautifulSoup
2 import urllib.request, urllib.error, urllib.parse
3 import re
4 import pandas as pd
5 from urllib.parse import urljoin
6
7
8
9
10
11
12 global all_books
13 all_books = []
14
15 def parse_book(url):
16     #url = 'https://books.toscrape.com/catalogue/tipping-the-velvet_999/index.html'
17
18     str = urllib.request.urlopen(url).read().decode()
19     str = str.strip()
20     soup = BeautifulSoup(str, 'html.parser')
21
22     #name = soup.find('body').find('div',class_ = 'container-fluid page').find('div',class_ = 'page_inner').find('div',class_ = 'content').find('h1').text.strip()
23     name = soup.find('div',class_ = 'content').find('h1').text.strip()
24
25     stock = soup.find('p',class_ = 'instock availability').text.strip()
26     stock = re.search(r'[0-9]+',stock).group()
27
28     rating = soup.find('p',class_='star-rating')['class'][1]
29     rnum = {
30         "One": 1, "Two": 2, "Three": 3, "Four": 4, "Five": 5
31     }
32     rating = rnum.get(rating,0)
33
34     genre = soup.find('ul',class_='breadcrumb').find('a',href = re.compile(r'category/books/')).text
35
36     p_items = dict()
37     p_items['Name']=name
38     p_items['Stock']=stock
39     p_items['Rating']=rating
40     p_items['Genre']=genre
41     t_items = [ 'UPC', 'Product Type', 'Price (incl. tax)', 'Price (excl. tax)', 'Tax', 'Number of reviews' ]
42     def tags(tag):
43         th_tag =soup.find('th',string=tag)
44
45         if (th_tag):
46             td =th_tag.find_next_sibling('td').text.strip()
47             else :
```

```
15 def parse_book(url):
16     p_items['Genre']=genre
17     t_items = ['UPC','Product Type','Price (incl. tax)','Price (excl. tax)','Tax','Number of reviews']
18
19 def tags(tag):
20     th_tag =soup.find('th',string=tag)
21
22     if (th_tag):
23         td =th_tag.find_next_sibling('td').text.strip()
24     else :
25         td = None
26     return td
27
28
29 for item in t_items:
30     p_items[item] = tags(item)
31 print(p_items)
32 all_books.append(p_items)
33 print(len(all_books))
34
35
36 def home_page(url):
37
38     #url = 'https://books.toscrape.com/catalogue/category/books/travel_2/index.html'
39
40     str = urllib.request.urlopen(url).read().decode()
41     str = str.strip()
42     soup = BeautifulSoup(str, 'html.parser')
43
44     genres = soup.find('ul',class_="nav nav-list").find_all('a')
45
46     links = list()
47     for genre in genres:
48         link = 'https://books.toscrape.com/' + genre['href']
49         links.append(link)
50
51     links = links[1:] # Removing the home link
52
53     for link in links:
54         book_link(link)
55
56
57
58
59 def book_link(url):
60     while True:
61         str_ = urllib.request.urlopen(url).read().decode().strip()
62         soup = BeautifulSoup(str_, 'html.parser')
63
64         # Find all book links on the current page
65         book_tags = soup.select('h3 a')
66         for tag in book_tags:
67             book_url = urljoin(url, tag['href'])
68             try:
69                 parse_book(book_url)
70             except:
71                 print(book_url, 'Not found')
72
73         # Check for next page
74         next_btn = soup.find('li', class_='next')
75         if next_btn:
76             next_link = next_btn.find('a')['href']
77             url = urljoin(url, next_link)
78         else:
79             break
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
```

Name	Stock	Rating	Genre	Price
The Perfect Play	4	3	Romance	£59.99
Last One Home	5	3	Fiction	£59.98
Civilization and Its Discontents	3	2	Psychology	£59.95
The Barefoot Contessa Cookbook	6	5	Food and Drink	£59.92
The Diary of a Young Girl	12	3	Nonfiction	£59.90
The Bone Hunters	4	3	Thriller	£59.71
Thomas Jefferson and the Tripoli Pirates	15	1	History	£59.64
Boar Island	12	3	Mystery	£59.48
The Improbability of Love	11	1	Fiction	£59.45
The Man Who Mistook His Wife for a Hat and Other Clinical Tales	6	4	Nonfiction	£59.45
The Gray Rhino	12	4	Add a comment	£59.15
Life Without a Recipe	3	5	Autobiography	£59.04
Listen to Me	1	3	Romance	£58.99
Unlimited Intuition Now	9	4	Default	£58.87
Approval Junkie	5	5	Autobiography	£58.81
Hamilton	14	3	Nonfiction	£58.79
Myriad	1	4	Fantasy	£58.75
The Rose & the Dagger	3	4	Fantasy	£58.64
Candide	4	3	Classics	£58.63
Alight	5	4	Default	£58.59
Catherine the Great	6	4	History	£58.55
Miller's Valley	8	2	Fiction	£58.54
Shameless	1	3	New Adult	£58.35
Silence in the Dark	8	3	Suspense	£58.33
How to Speak Golf	12	5	Default	£58.32
Aristotle and Dante Discover the Secrets of the Universe	14	4	Young Adult	£58.14
The Death of Humanity	16	4	Philosophy	£58.11
Unstuffed	3	1	Nonfiction	£58.09



Sum of Stock	Column Labels					
Row Labels	1	2	3	4	5	Grand Total
Academic		5				5
Add a comment	99	117	107	121	72	516
Adult Fiction					3	3
Art		7	22	42	1	72
Autobiography	22	21		3	11	57
Biography	7	26	10			43
Business	27	18	29	19	40	133
Childrens	62	41	79	27	20	229
Christian	7	14			16	37
Christian Fiction		7		18	30	55
Classics	7	20	15	13		55
Contemporary	19				14	33
Crime	15					15
Cultural	15					15
Default	260	341	266	238	240	1345
Erotica					15	15
Fantasy	82	61	82	60	87	372
Fiction	119	54	168	102	145	588
Food and Drink	90	62	71	32	64	319
Health		14	16		19	49
Historical	6				14	20
Historical Fiction	48	29	39	32	46	194
History	47	33	50	11	40	181
Horror	44	37	26	12	17	136
Humor	21	3	12	14	34	84



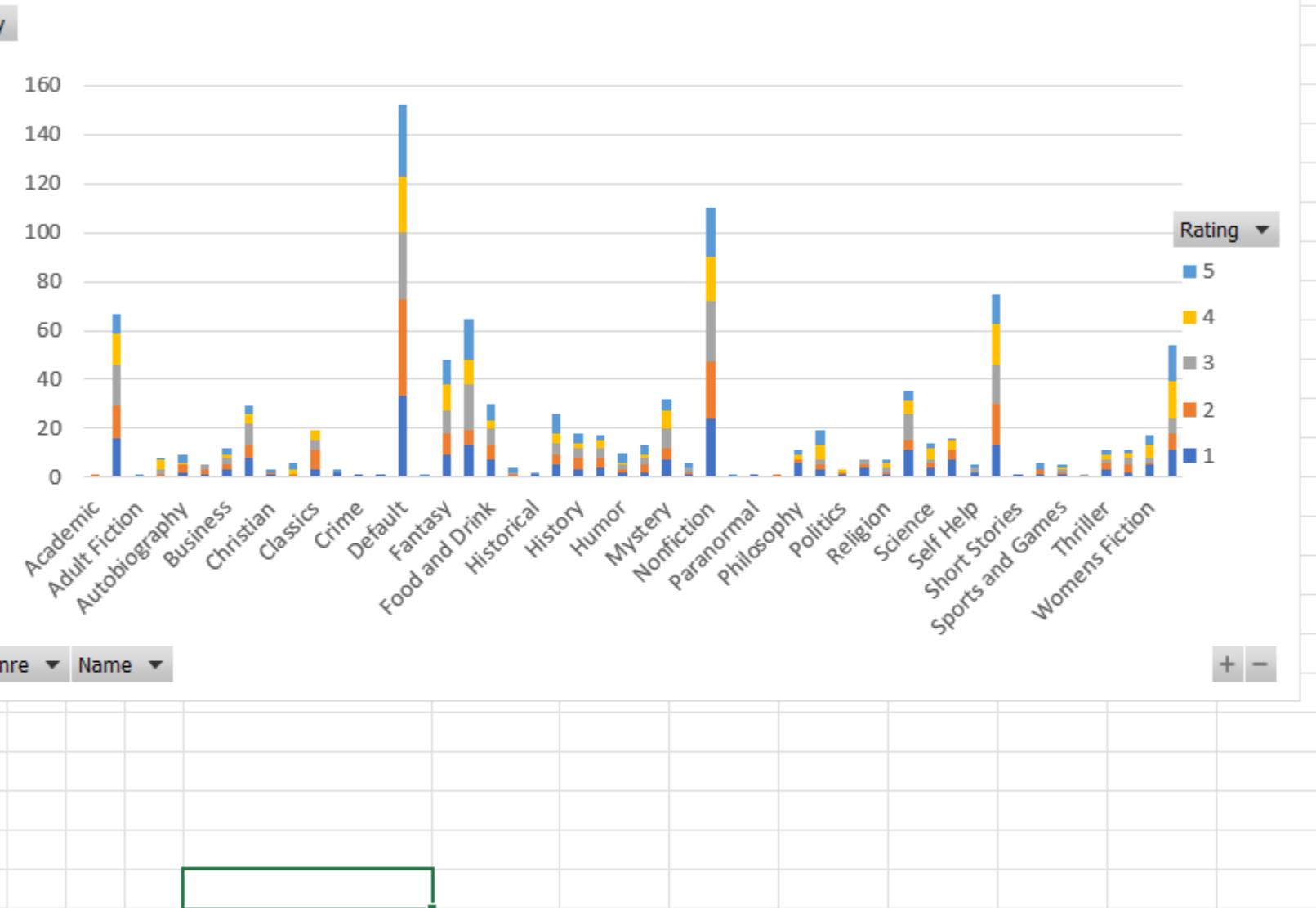
Rating

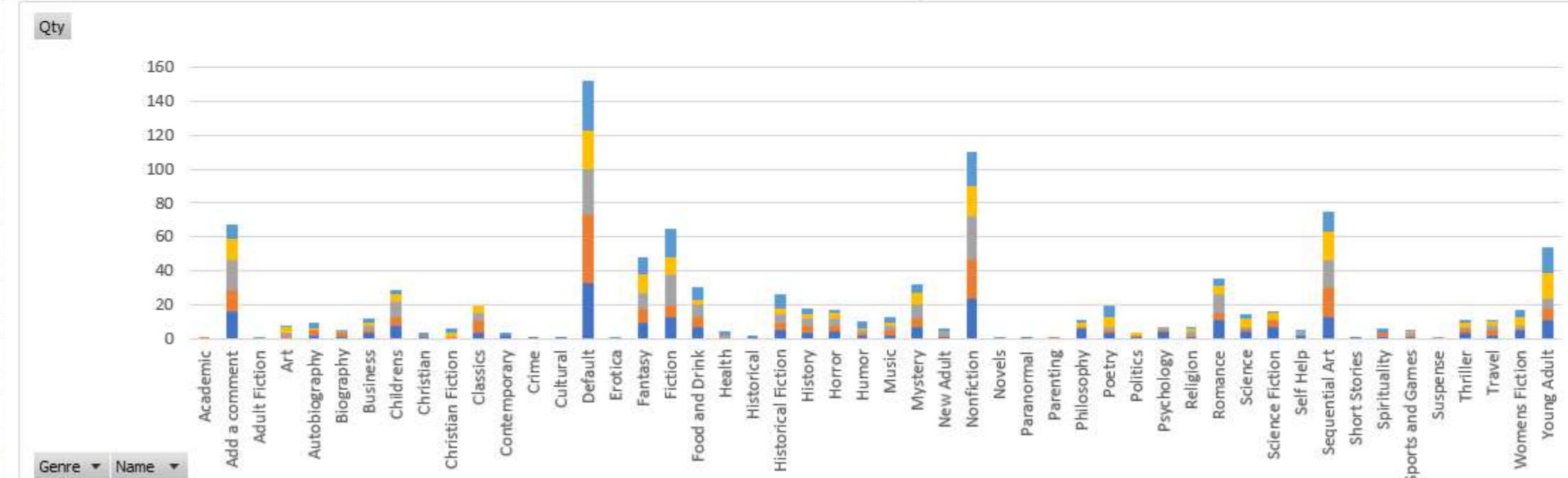
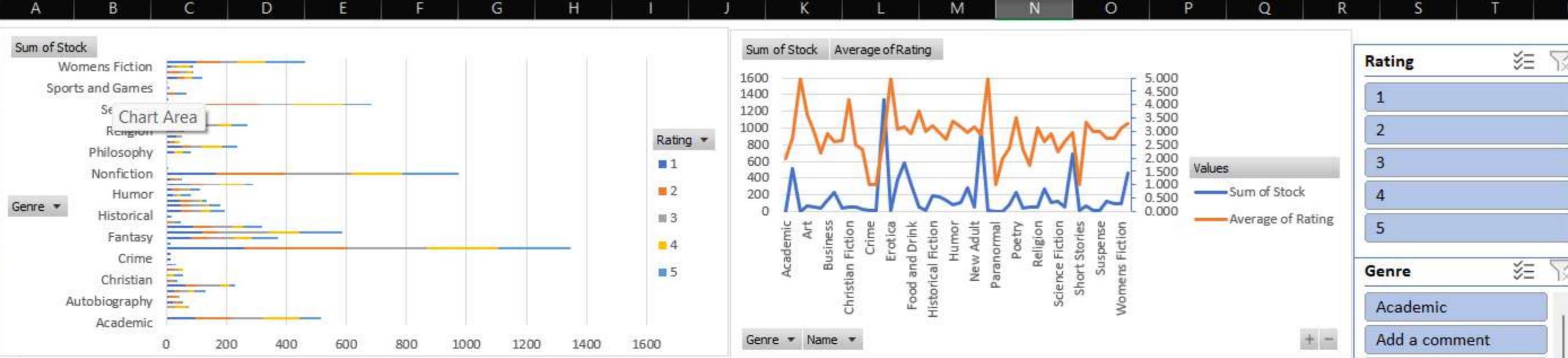
1
2
3
4
5

Genre

Academic
Add a comment
Adult Fiction
Art
Autobiography
Biography
Business
Childrens

Qty	Column Labels					
Row Labels	1	2	3	4	5	Grand Total
Academic	1					1
Add a comment	16	13	17	13	8	67
Adult Fiction				1		1
Art		1	2	4	1	8
Autobiography	2	3	1	3		9
Biography	1	2	2			5
Business	3	2	3	1	3	12
Childrens	8	5	9	4	3	29
Christian	1	1			1	3
Christian Fiction		1		2	3	6
Classics	3	8	4	4		19
Contemporary	2			1		3
Crime	1					1
Cultural	1					1
Default	33	40	27	23	29	152
Erotica				1		1
Fantasy	9	9	9	11	10	48
Fiction	13	6	19	10	17	65
Food and Drink	7	6	7	3	7	30
Health	1	1		2		4
Historical	1			1		2
Historical Fiction	5	4	5	4	8	26
History	3	5	4	2	4	18
Horror	4	4	4	3	2	17
Humor	2	1	2	1	4	10
Mystery	6	6	6	1	1	19
Nonfiction						
Paranormal						
Philosophy						
Politics						
Religion						
Science						
Sports and Games						
Thriller						
Womens Fiction						





➤ Project - 01 (DA-01) :-



to scrape "books.firebaseio.com", prepare an excel sheet & derive conclusions :-

- i) Price v/s rating
- ii) Distribution of Star Ratings
- iii) Stock Availability

• features :-

- Extracts all books across all categories.
- Collects detailed meta data :-

- Name
- Stock
- Rating
- UPC
- Product type
- Price (exc. tax)
- Price (inc. tax)
- Tax
- No. of reviews
- Genre

- Handles pagination & save it in CSV.

• How it works ??

1. Home Page Parsing :-

Takes home page url
 & give links for
 all genre's page.

`def home_page(url) :`

Created S : `S = urllib.request.urlopen(url).read().decode()`
 soup [`S = S.strip()`
 of given `soup = BeautifulSoup(S, 'html.parser')`
 (url).

`genres = soup.find('ul', class_="nav nav-list")
 .find_all('a')`

Searched targeted class & tag.

`links = list()`

for genre in genres :

`link = 'https://books.toscrape.com/' + genre['href']
 links.append(link)`

In genre['href'] we only get the relative address, so we add if if rest.

`links = links[1:] # Removed homepage link`

for link in links :

`book-link(link)`

things to import ??

```
from bs4 import BeautifulSoup  
import urllib.request, urllib.error, urllib.parse  
import re  
import pandas as pd  
from urllib.parse import urljoin
```

created a all-books dict.

global all_books
all_books = {}

2. Genre Crawling :-

- Iterates through each genre & fetches all book listing pages.
- To Pagination (if next page appears moves to next).

def book_list(url) :

 while True :

 # created soup from url.

find all books ↓

book_tags = soup.select('h3>a')

for tag in book_tags :

 book_url = urljoin(url, tag['href'])

used a func? rather manual join & merge.

try :

 Parse_book(book_url)

except :

 print(book_url, 'Not found!')

for next page

next_btn = soup.find('li', class_="next")

if next_btn :

 next_link = next_btn.find('a')['href']

 url = urljoin(url, next_link)

else :

 break

3. Parsing a book :-

→ Visit each book page & extract the details.

def parse_book(curl):

soup created from curl

Extracted all parameters :-

name = soup.find('div', class_ = 'content').find('h1')
 • text.strip()

stock = soup.find('p', class_ = 'inStock availability')
 • text.strip()

stock = re.search(r'[0-9]+', stock).group
 # to extract qty. in numbers.

rating = soup.find('p', class_ = 'star rating')[1].

num = { "One": 1, "Two": 2, "Three": 3,
 "Four": 4, "Five": 5
 } # text → Num

rating = num.get(rating, 0)

- # Created a dic. to take output extracted values (P-items)
- # Created dic. to iterate in single table (t-items)

P-items = dict()

P-items ['Name'] = name

P-items ['Stock'] = stock

P-items ['Rating'] = rating

P-items ['Genre'] = genre

t-items = ['UPC', 'Product Type', 'Price (inc. tax)',
 'Price (exc. tax)', 'Tax', 'Number
 of Reviews']

created func? (tags) to iterate each item.

def tags(tag):

th-tag = soup.find('th', string=tag)

↑ find inside with respected
 th-table. names.

if (th-tag):

td = th-tag.final-next-sibling('td')
 .text.strip()

else:

td = None

return td

for item in t-items:
 P-items[item] = tags(item)

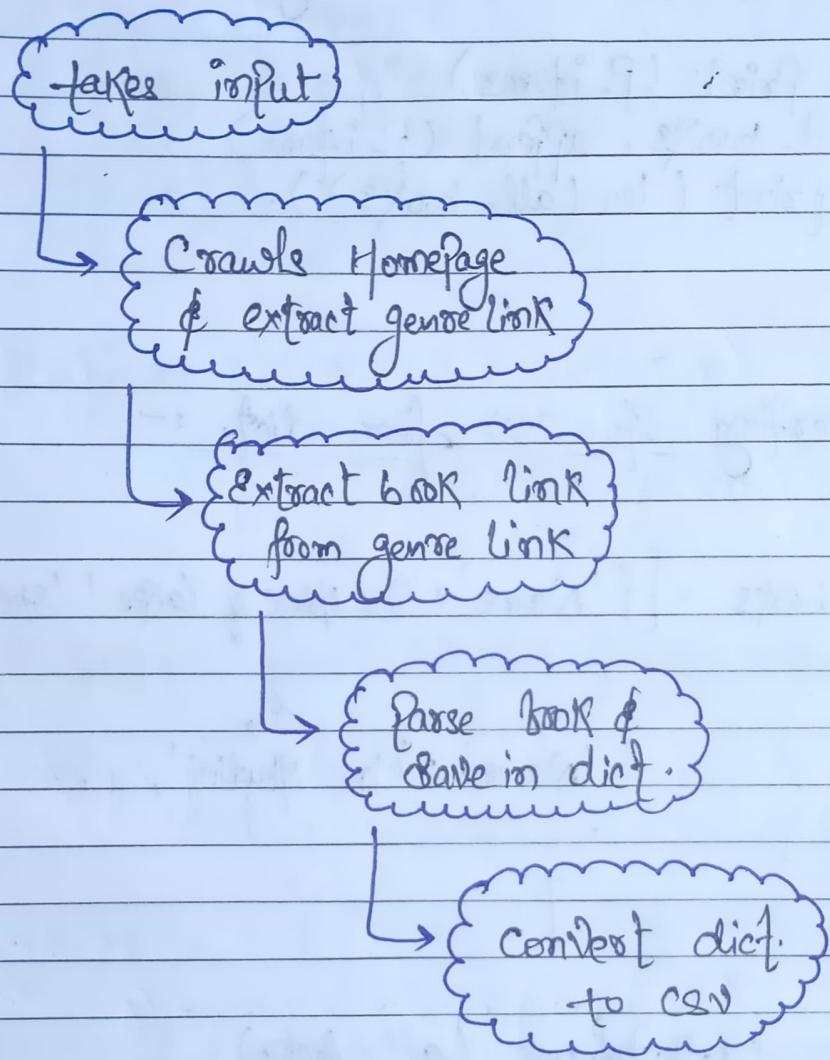
```
# Print (P-items)
all-books.append (P-items)
# Print (len (all-books)).
```

4. Converting to CSV from dict. :-

```
# all-books = [{ 'Name': 'The House of Cards', 'Stock': 15
                },
                { 'Name': 'Mrs. Maudlin',
                }]
```

```
df = pd.DataFrame (all-books)
df.to_csv ("books-data.csv"), index = False )
```

• Program workflow :-



- Methods () used :-

- Soup.text.strip()

:- remove tag & prints the inside text.

- title = soup.find('head').find('title')

:- iterate over tags to find req. element.

:- We do it many times (nesting)

Eg:-

```
name = soup.find('body').find('div', class_='container fluid page').find('div', class_ = 'page-inner').find(...).find('h1').text.strip()
```

- re.search()

:- returns index & other details of matched content.

- `re·Search().group()`
:- returns only matched string.
- `rating = soup.get('rating', 0)`
:- To extract values from Key inside dict.
- `soup.find(tag_name, attribute filter)`

Eg:-

`N-review = soup.find('a', href = re.compile(r'reviews/'))`

↳ find 'a'

↳ inside match attribute

- Prettify()

- :- to beautify HTML codes.

- :- but we don't use methods in it.

- to Parse HTML Comment

- <!-- text --> comment

```
from bs4 import BeautifulSoup, comment
```

→ Search all Comments :-

Using

```
soup.findAll(string = lambda text:  
isinstance(text, comment))
```

- :- returns normal HTML text.

- :- Parse it as usual.

- `upc = soup.find('h1', attrs = {"UPC": "})`
- ↑
tag name. ↑ attribute filter

In UPC,
 ↳ it doesn't have any class or
 ↳ tag.

So,
 we navigate by value.

- `link = genre['href']`

`genre = `
 Books
``

↳ To iterate links inside tag.

Point (link)

↳html

↳ Clean website.

> Cleaning imported CSV with Pandas :-

Issues :-

- ① Big Names
- ② Missing Values
- ③ Extra Spaces
- ④ Duplicates
- ⑤ type cast
- ⑥ clean price, remove symbol
- ⑦ Remove unnecessary data
 - ↳ UPC :- no use
 - ↳ N.o. of reviews :- 0
 - ↳ Page :- 0
 - ↳ Product type :- Books.

Price & Price :- contains same value {:- 250}

{inc. tax} {exc. tax}

merged to "Price".

codes :-

```
import Pandas as pd
import re

df = pd.read_csv('... .csv') # import
df.head() # Analyze
df.info()
```

```
df['Name'] = df['Name'].str.replace(':', '').str[0].str.strip()
df['Name'] = df['Name'].str.replace(r'[^a-zA-Z0-9]', '',
                                     regex=True).str.strip()
```

trimmed name after : (.) inside this.

```
df[df.duplicated()]
# no duplicates found
```

remove irrelevant columns.

```
df = df.drop('Tax', axis=1)
df = df.drop('Number of Reviews', axis=1)
df = df.drop('Price (excl. Tax)', axis=1)
df = df.drop('Product Type', axis=1)
```

rename 'Price (incl. Tax)' → 'Price'.

```
df = df.rename(columns={'Price (incl. Tax)': 'Price'})
df['Price'] = df['Price'].str[1:] # remove symbol.
```

```
df.to_csv('new-name.csv', index=False)
```

save

Spiral

➤ Excel Analyzing & Dashboarding :-



→ Created 4 Sheet & a dashboard :-

- (i) Cleaned books chart
- (ii) Stock by rating
- (iii) Rating by Stock
- (iv) Rating by Qty.

Dashboard

→ added Slicers & graphs.