

INDENG 243 Project - Module 2

GoodReads: Book Analytics & Recommendation System

- @Group: 18

Agenda

- Model Data Preparation
- Recommendation models
 - Popularity Model (Baseline)
 - Content-based Filtering
 - Collaborative Filtering
 - Hybrid Model
- NLP Model
 - Tag filtering system
 - Extractive summarization
- Summary & Future plans

Part 1) Model data preparation

```
In [1]: import pandas as pd
import numpy as np

df_original = pd.read_csv('./Datasets/Book_info_original.csv')
```

1.1 Data Revise

In previous exploratory data analysis, we observed some potential flaws in original datasets:

1.1 Insufficient *manga* books

- Books under genre *manga* seems being especially insufficient in data (only 20 books) compared to other genres (100~300 books)

```
In [10]: len(df_original[df_original['Genre'] == 'manga'])

Out[10]: 20
```

1.2 Hidden duplicated books

- Interestingly, some books are implicitly duplicated in original data: they only differ in `Genre` but all these records refer to a same book
- The reason is that in our previous web scraping, for simplicity we keep the pre-defined 40 genres as the main `Genre` attribute for each book, but sometimes a single book will be listed under many different genres, resulting in *duplicated* records
 - Eg. See the following *Harry Potter* and *The Chamber of Secrets*, it actually was listed under fantasy, fiction, children, etc
- However, these implicit duplicates are not detected by codes directly using ***duplicated()***

```
In [12]: # if directly using .duplicated() to detect duplicates
df_original.duplicated().sum()

Out[12]: 0
```

```
In [15]: df_original[df_original['UId'] == 15881].head(4)
```

	UId	Title	Author	Genre	Sub_Genres	Rating	Publish_Date	Page_Num	Award	Description	Author_Desc	Book_Au
	974	15881	Harry Potter and the Chamber of Secrets	J.K. Rowling	children-Young Adult, 'Magic', 'Childrens...	4.43	July 2, 1998	341	[Mythopoeic Fantasy Award, 'British Book Awa...	Ever since Harry Potter had come home for the ...	See also: Robert GalbraithAlthough she writes...	
	2509	15881	Harry Potter and the Chamber of Secrets	J.K. Rowling	ebooks	4.43	July 2, 1998	341	[Mythopoeic Fantasy Award, 'British Book Awa...	Ever since Harry Potter had come home for the ...	See also: Robert GalbraithAlthough she writes...	
	2671	15881	Harry Potter and the Chamber of Secrets	J.K. Rowling	fantasy	4.43	July 2, 1998	341	[Mythopoeic Fantasy Award, 'British Book Awa...	Ever since Harry Potter had come home for the ...	See also: Robert GalbraithAlthough she writes...	
	2869	15881	Harry Potter and the Chamber of Secrets	J.K. Rowling	fiction	4.43	July 2, 1998	341	[Mythopoeic Fantasy Award, 'British Book Awa...	Harry Potter had come home for the ...	See also: Robert GalbraithAlthough she writes...	

Subsequently this observatin gives a rise to other fundamental problems:

- 1) *Which of these duplicated records should we keep?* (i.e. Which genre you think these books should be?)
 - It's hard to answer because a book's content might interact with many different genres and randomly keeping one record is unreasonable
 - So we decide to obtain each book's genre list by re-scrape it under book's information page, rather than genre classification given by website
- 2) If we drop the duplicates, the number of unique books are greatly reduced (8032 --> 5545)

```
In [136]: print("The original total book number is: ", len(df_original))
print("The unique book number is: ", len(df_original.drop_duplicates(subset=['UId'], 'Title'))))

The original total book number is: 8032
The unique book number is: 5545
```

1.3 Biased Reviews

- Review dataset might be biased: more positive reviews than negative ones (feedback from Module 1)
- Currently **# Negative review (Rating: 1~3) : # Positive review (Rating: 4~5) = 1 : 4** (roughly), unbalanced

To solve the above problems, we decide to re-scrape additional book data, which is in nature the best and straightforward solution

The revised dataset has the following improvements targeting previous problems:

- More *manga* books: 20 --> 140
- Re-scrape the book genre list via re-detailed code in `UId_Genre.py`
- Re-scrape more unique books: 5545 --> 7846 (`Genre_Book_Add.py`)
- Correspondingly more book reviews: 80461 --> 115033 and we purposely select more negative reviews to make dataset more balanced
 - **# Negative review (Rating: 1~3) : # Positive review (Rating: 4~5) = 1 : 2** (more balanced compared to before)

For data combination, cleaning and some pre-processing steps, since they are lengthy and not close to this module's topic, you can see detailed codes in *Data_Read_Cleaning.py* if you like. We will directly use the further cleaned datasets in following sections:

- **Book_info_cleaned.csv**
- **Book_reviews_cleaned.csv**
- **Book_stats** (folder for all book statistics, since 1 book has 4 x 170 entries, we keep them separately)

1.2 Datasets modification for model use

1.2.1 Book information dataset: df_info dataframe

```
In [16]: df_info = pd.read_csv('./Datasets/Basic_datasets/Book_info_cleaned.csv')
```

```
In [17]: df_info.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7776 entries, 0 to 7775
Data columns (total 19 columns):
 # Column Non-Null Count Dtype
---
 0 UId 7776 non-null int64
 1 Title 7776 non-null object
 2 Author 7776 non-null object
 3 Rating 7776 non-null float64
 4 Publish_Date 7776 non-null object
 5 Page_Num 7776 non-null int64
 6 Award 7776 non-null object
 7 Description 7761 non-null object
 8 Author_Desc 7116 non-null object
 9 Book_Authorized 7776 non-null int64
 10 Follower_Num 7776 non-null int64
 11 Review_Num 7776 non-null int64
 12 Rating_Num 7776 non-null int64
 13 Four_Star 7776 non-null int64
 14 Three_Star 7776 non-null int64
 15 One_Star 7776 non-null int64
 16 N_Likes 7776 non-null int64
 17 N_Comments 7776 non-null int64
dtypes: float64(1), int64(11), object(7)
memory usage: 1.1+ MB
```

Turn rating distribution into percentage

- Since different books have different volume of ratings number, it is more reasonable to look at the percentage of each rating category rather than the real rating number. So we first compute the rating percentages as follows:

```
In [18]: df_info['Five_star_percent'] = df_info['Five_Star'] / df_info['Rating_Num']
df_info['Four_star_percent'] = df_info['Four_Star'] / df_info['Rating_Num']
df_info['Three_star_percent'] = df_info['Three_Star'] / df_info['Rating_Num']
df_info['Two_star_percent'] = df_info['Two_Star'] / df_info['Rating_Num']
df_info['One_star_percent'] = df_info['One_Star'] / df_info['Rating_Num']

In [130]: df_info = df_info.drop(['Five_Star', 'Four_Star', 'Three_Star', 'Two_Star', 'One_Star'], axis=1)
```

Data pre-processing before using df_info dataframe

- Because the list-type data are stored as *strings*, we need to turn them back to the correct data type.

```
In [28]: from ast import literal_eval

df_info['Genres'] = df_info['Genres'].apply(lambda x: literal_eval(x) if "[]" in x else x)
df_info['Award'] = df_info['Award'].apply(lambda x: literal_eval(x) if "[]" in x else x)
```

Award - Ordinal Encoding

- We use the number of awards the book has obtained to represent the 'Award' attribute
- <https://analyticsindamag.com/a-complete-guide-to-categorical-data-encoding/>

```
In [21]: Award_Num = df_info['Award'].apply(lambda x: 0 if x == ['no'] else len(x))
df_info['Award_Num'] = Award_Num
```

We look at the book ratings' distribution

```
In [48]: print("The number of low rating books (Rating <= 4.0): ", len(df_info[df_info['Rating'] < 4.0]))
print("The number of high rating books (Rating > 4.0): ", len(df_info[df_info['Rating'] >= 4.0]))

The number of low rating books (Rating <= 4.0): 2992
The number of high rating books (Rating > 4.0): 4784
```

Book Statistics

The book statistics data reflects the *changes* (compared to previous day) in the following attributes over past 6 months:

- **date**: The date for these statistics
 - **added**: changes of times this book has been added to users' bookshelves from previous day
 - **ratings**: changes of times this book has been rated by users from previous day
 - **reviews**: changes of times this book is left reviews by users from previous day
 - **to-read**: changes of times this book has been marked as 'to-read' by users
- (Note: the first row in file is the sum of these statistics until the day we scraped the data)

We want to combine book statistics into its information dataframe, using `df_info`, however, each book only has one entry in information dataframe, while it has 4 x 170 records for its statistics. Directly combining them will make the dimensionality messy and hard to interpret, so we decide to use their statistical indicators for better measurement:

- **Median**: measures the central tendency of data distribution
- **IQIR (Interquartile Range)**: measures the dispersion level of data distribution (*more robust to scale problem compared to variance*)
 - https://en.wikipedia.org/wiki/Robust_measures_of_scale

```
In [36]: import os

path = './Partial_datasets/Book_stats/'
filenames = os.listdir(path)
```

```
In [68]: df_demo = pd.read_csv(path + filenames[0])
df_demo.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 172 entries, 0 to 171
Data columns (total 19 columns):
 # Column Non-Null Count Dtype
---
 0 date 172 non-null object
 1 added 172 non-null int64
 2 ratings 172 non-null int64
 3 reviews 172 non-null int64
 4 to-read 172 non-null int64
dtypes: int64(4), object(1)
memory usage: 6.8+ MB
```

Combine book statistics and compute statistical indicators simultaneously

```
In [38]: df_info_stats = pd.DataFrame()

for file in filenames:
    uid = int(file.split('_stats')[0])
    stats_df = pd.read_csv(path + file)

    # if the first line is the sum of all following statistics, so separately processing it
    data = stats_df.iloc[1:, 1:].astype('float64').values
    data_df = pd.DataFrame(data, columns=['UId', 'AddedShelf', 'ToRead'])

    # compute median & IQR
    stats = stats_df.iloc[1:, 1:]
    IQR_df = (stats.quantile(q=0.75) - stats.quantile(q=0.25)).add_suffix('_IQR')
    median_df = stats.median().add_suffix('_median')

    stats_combine = pd.concat([data_df, median_df, IQR_df], axis=1)
    df_info_stats = pd.concat([df_info_stats, stats_combine], ignore_index=True)

In [39]: df_info_stats[['UId', 'AddedShelf', 'ToRead']] = df_info_stats[['UId', 'AddedShelf', 'ToRead']].astype('int64')
df_info_stats['UId'] = df_info_stats['UId'].astype('int64')
```

```
Out[39]:
```

	UId	AddedShelf	ToRead	added_median	ratings_median	reviews_median	to-read_median	added_IQR	ratings_IQR	reviews_IQR	rean
0	10006486	13107	14731	7.0	0.0	0.0	4.0	1.00	2.00	0.0	
1	10006486	2627	6709	7.0	0.0	0.0	4.0	1.00	0.00	0.0	
2	10005151	124063	926493	24.0	7.0	1.0	7.0	8.00	7.00	0.0	
3	10008056	36918	6637	12.0	4.0	0.0	2.0	6.00	4.75	0.0	
4	10003977	27193	11902	22.0	6.0	0.0	9.0	17.75	5.00	1.0	
...
7771	99713	22569	8889	3.0	1.0	0.0	2.0	3.00	1.00	0.0	
7772	9984	244540	121383	101.0	28.0	1.0	59.0	39.00	12.00	2.0	
7773	9994	192255	98182	58.0	10.0	1.0	35.0	18.00	7.00	1.0	
7774	99955	35739	18080	5.0	1.0	0.0	3.0	2.50	1.00	0.0	
7775	9999107	85198	32656	19.0	5.0	0.0	7.0	7.75	5.00	1.0	

7776 rows x 11 columns

Combine statistics data with original df_info dataframe

```
In [41]: df_info = df_info.merge(df_info_stats, on='UId')
```

```
In [42]: # Store as csv for further use
df_info.to_csv('./Datasets/Basic_datasets/Book_info_model.csv', index=False)
```

1.2.2 Book reviews dataset: df_reviews dataframe

```
In [117]: df_reviews = pd.read_csv('./Datasets/Basic_datasets/Book_reviews_cleaned.csv')
df_reviews.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 114034 entries, 0 to 114033
Data columns (total 10 columns):
 # Column Non-Null Count Dtype
---
 0 UId 114034 non-null int64
 1 Title 114034 non-null object
 2 Reviewer 114034 non-null object
 3 N_Review 114034 non-null int64
 4 N_Follower 114034 non-null int64
 5 Review_Rating 114034 non-null int64
 6 Review_Date 114034 non-null object
 7 Content 114034 non-null object
 8 N_Likes 114034 non-null int64
 9 N_Comments 114034 non-null int64
dtypes: int64(6), object(4)
memory usage: 8.7+ MB
```

Generally we keep the same ratio of negative : positive ratings as book information datasets (around 1 : 2)

- ```
In [118]: print("The number of negative reviews (Rating: 1~3): ", len(df_reviews[df_reviews['Review_Rating'] < 4]))
print("The number of positive reviews (Rating: 4~5): ", len(df_reviews[df_reviews['Review_Rating'] >= 4]))

The number of negative reviews (Rating: 1~3): 32003
The number of positive reviews (Rating: 4~5): 78831
```

We want to examine our *unique* users are the same in our further recommendation models

- Here *unique* means these users not only share the same name but also the same number of reviews (`Reviewer' + N_Review`)
- Because GoodReads allow users to use the same username to register an account, we would like to use the `N_Review` (# of reviews written to help identification)

```
Out[119]: df_reviews.groupby(by=['Reviewer', 'N_Review']).count()

Reviewer N_Review
Amanda 31 1 1 1 1 1 1 1 1
Luna 265 4 4 4 4 4 4 4 4
(than) Littlebookcove 112 1 1 1 1 1 1 1 1
The Polybrary 313 1 1 1 1 1 1 1 1
A. 1163 11 11 11 11 11 11 11 11
...
점스마 165 1 1 1 1 1 1 1 1
여리고 167 1 1 1 1 1 1 1 1
D r e y 71 6 6 6 6 6 6 6 6
P r e k k e 845 17 17 17 17 17 17 17 17
e l l i e 182 1 1 1 1 1 1 1 1
```

41271 rows x 8 columns

But the above result of unique users is actually inaccurate, please see the following example:

- These two users with the same name 'Luna' but different number of reviews are actually the **same** user (we checked on website)
- But why we see this situation happens? The reasons are the following:
  - 1) GoodReads website has some historical problem about user information updates, i.e. the user might have more review numbers than some historical timestamps, but the website did not update these information cascade;
  - 2) Our web scraping takes some time, the user's review number might have changed in this time interval, making them *different*.

P.S. You might notice that they have different `N_Follower` numbers, that is caused by web scraping error, we will solve it in the following part

```
In [121]: df_reviews[df_reviews['Reviewer'] == 'Luna']

Out[121]:
```

| UId    | Title  | Reviewer                                   | N_Review | N_Follower | Review_Rating | Review_Date | Content         | N_Likes                                           | N_Comments |   |
|--------|--------|--------------------------------------------|----------|------------|---------------|-------------|-----------------|---------------------------------------------------|------------|---|
| 36645  | 47944  | Empress of the World                       | Luna     | 638        | 43            | 4           | August 21, 2011 | There's not one thing that makes this book goo... | 3          | 2 |
| 110764 | 112118 | The Cold Heaven Seven Seasons in Greenland | Luna     | 645        | 4             | 4           | April 6, 2011   | Oh Greenland. Sometimes I'll get there and sa...  | 5          | 0 |

#### First we modify the user names for better matching

- Some Reviewers actually are the same person, but they are actually not the same names, how can we distinguish them? We will discuss this problem shortly after this. We first just modify the reviewer names for better format
- Actually the following user with name 'Luna' is an example, the 1st and 5th are the same user, while the other 4 are a different user

```
In [128]: for index, row in df_reviews.iloc[[36645, 59944, 60198, 60782, 110764, 110967], 2].iteritems():
 print(index, row)

(36645, 'Luna')
(59944, 'Luna')
(60198, 'Luna')
(60782, 'Luna')
(110764, 'Luna')
(110967, 'Luna')
```

```
In [129]: df_reviews.iloc[[36645, 59944, 60198, 60782, 110764, 110967], :2]

Out[129]:
```

| UId    | Title  | Reviewer                                         | N_Review | N_Follower | Review_Rating | Review_Date | Content            | N_Likes                                           | N_Comments |   |
|--------|--------|--------------------------------------------------|----------|------------|---------------|-------------|--------------------|---------------------------------------------------|------------|---|
| 36645  | 47944  | Empress of the World                             | Luna     | 638        | 43            | 4           | August 21, 2011    | There's not one thing that makes this book goo... | 3          | 2 |
| 59944  | 203310 | Poetry Language Thought                          | Luna     | 265        | 16            | 4           | September 22, 2014 | This was a refreshing read. What I really like... | 2          | 0 |
| 60198  | 232743 | The Archeology of Knowledge and the Discourse... | Luna     | 265        | 16            | 2           | November 27, 2016  | This is the sort of book that you feel that I...  | 4          | 0 |
| 60782  | 252648 | The Postmodern Condition A Report on Knowledge   | Luna     | 265        | 16            | 5           | February 14, 2016  | Postmodernism. What now? After the fall of the... | 2          | 0 |
| 110764 | 112118 | This Cold Heaven Seven Seasons in Greenland      | Luna     | 645        | 4             | 4           | April 6, 2011      | Oh Greenland. Sometimes I'll get there and sa...  | 5          | 0 |
| 110967 | 117160 | Three Dialogues Between Hylas and Philonous      | Luna     | 265        | 1             | 5           | September 16, 2021 | Basically the material world does not exist pe... | 1          | 0 |

```
In [130]: from string import punctuation

def remove_punctuation(document):
 document = document.replace(' ', '') # remove empty space
 no_punct = ''.join(character for character in document if character not in punctuation)

 if no_punct == '':
 return document
 else:
 return no_punct
```

```
In [131]: df_reviews['Reviewer'] = df_reviews['Reviewer'].apply(remove_punctuation)
```

```
In [132]: # reviewer names
names = df_reviews['Reviewer'].unique()
print(len(names))

23547
```

The number of 'unique' (name + review number) users are not much different from before

- To some extent, we would say our operation of modifying name format did not make too much negative impacts on data integrity

```
In [133]: print("The number of unique users are: ", len(df_reviews.groupby(by=['Reviewer', 'N_Review']).count()))

The number of unique users are: 41255
```

To tackle this problem, our approach is to combine reviewers' entries with the same name & close review numbers

- We assume that the user with the same name and close number of reviews are the same user.
- From deep observation and several trials on original data, we set 100 as the limit for the definition of close reviews (i.e. absolute difference is smaller or equal to 100), and we will check the quality of this operation afterwards.

```
In [382]: # we use the copy to avoid directly operating on original data
df_reviews_copy = df_reviews.copy()

In []: for name in names:
 i = 1
 temp_group = []
 temp_non_group = []
 review_group = [0]
 non_review_group = [0]

 flag = False
 while (temp_group != review_group) & (temp_non_group != non_review_group):
 print("Current user name: ", len(temp_group), len(temp_group))
 temp_group = [0]
 temp_non_group = [0]
 temp_review_group = []
 temp_non_review_group = []

 if flag:
 break

 same_name_entries = df_reviews_copy[df_reviews_copy['Reviewer'] == name]
 if len(same_name_entries) > 1:
 index = list(same_name_entries.groupby(by=['N_Review']).count())['UId']
 if len(temp_group) > 1:
 flag = True
 continue

 # find more than one entry user's review number
 if i > 1:
 review_group = list(temp[temp > 1].index.values)
 non_review_group = list(temp[temp == 1].index.values)
 else:
 review_group = temp_group[:]
 non_review_group = temp_non_group[:]

 # if exists aggregated entries
 remove_group = []
 review_group = []

 if review_group:
 for group v.s. independent sample
 for review_num in review_group:
 if np.abs(review_num - review_group) <= 100:
 same_name_entries = df_reviews_copy[df_reviews_copy['Reviewer'] == name]
 if review_num_1 < review_num_2:
 index = list(same_name_entries[same_name_entries['N_Review'] == \
 review_num_1].index.values)
 df_reviews_copy.loc[index, 'N_Review'] = review_num_2 # change values
 remove_group.append(index) # drop already allocated ones
 else:
 index = list(same_name_entries[same_name_entries['N_Review'] == \
 review_num_2].index.values)
 df_reviews_copy.loc[index, 'N_Review'] = review_num_1 # change values
 remove_group.append(review_num)

 print("-----")
 review_group = [i for i in review_group if i not in set(remove_group)]
 print(review_group)
 print(non_review_group)

 # group v.s. independent sample
 for non_group in non_review_group:
 if np.abs(review_num - non_group) <= 100:
 same_name_entries = df_reviews_copy[df_reviews_copy['Reviewer'] == name]
 if review_num > non_group:
 index = list(same_name_entries[same_name_entries['N_Review'] == \
 review_num].index.values)
 df_reviews_copy.loc[index, 'N_Review'] = non_group_2 # change values
 remove_group.append(index)
 else:
 index = list(same_name_entries[same_name_entries['N_Review'] == \
 review_num_1].index.values)
 df_reviews_copy.loc[index, 'N_Review'] = non_group_1 # change values
 remove_group.append(index)

 print("-----")
 non_review_group = [i for i in non_review_group if i not in set(remove_group)]
 print(non_review_group)
 print(non_review_group)

 # independent sample v.s. independent sample
 if not non_review_group:
 flag = True
 continue
 else:
 for non_group_1 in non_review_group:
 for non_group_2 in non_review_group:
 if np.abs(non_group_1 - non_group_2) <= 50:
 same_name_entries = df_reviews_copy[df_reviews_copy['Reviewer'] == name]
 if non_group_1 < non_group_2:
 index = list(same_name_entries[same_name_entries['N_Review'] == \
 review_num_1].index.values)
 df_reviews_copy.loc[index, 'N_Review'] = non_group_2 # change values
 remove_group.append(index)
 else:
 index = list(same_name_entries[same_name_entries['N_Review'] == \
 review_num_2].index.values)
 df_reviews_copy.loc[index, 'N_Review'] = non_group_1 # change values
 remove_group.append(index)

 print("-----")
 non_review_group = [i for i in non_review_group if i not in set(remove_group)]
 print(non_review_group)
 print(non_review_group)

 i += 1
```

```
In [135]: print("The number of unique users after combination are: ", len(df_reviews_copy.groupby(by=['Reviewer', 'N_Review']).count()))

The number of unique users after combination are: 27581
```

#### Solve the error caused in web scraping

- Minor mistake in regular expressions causing the `N_Follower` data is wrong: only matched the first digit
- We use the previous combination result, taking the highest `N_Follower` number of the same user to replace all entries

```
In [29]: for name in names:
 same_name_entries = df_reviews_copy[df_reviews_copy['Reviewer'] == name]
 # only operate on usernames with more than 1 entries
 if len(same_name_entries) > 1:
 temp = same_name_entries.groupby(by=['N_Review']).count()['UId']
 # if they already all shares the same N_Review, jump to next one
 if len(temp) == 1:
 continue
 review_group = list(temp[temp > 1].index.values)
 for review_num in review_group:
 same_review_df = same_name_entries[same_name_entries['N_Review'] == review_num]
 # find the maximum N_Follower as the most updated data to replace all other for same user
 max_follower_num = [same_review_df['N_Follower'].max()]
 index = list(same_name_entries[same_name_entries['N_Review'] == review_num].index.values)
 dup_follower_num = [i for i in max_follower_num for _ in range(len(index))] # duplicate for len
 df_reviews_copy.loc[index, 'N_Follower'] = dup_follower_num
```

#### Assign a unique UserID for each different user to better distinguish them

```
In [222]: df_reviews_copy.loc[:, 'UserID'] = df_reviews_copy.groupby(['Reviewer', 'N_Review']).nngroup()

In []: # adjust the column display
cols = df_reviews_copy.columns.to_list()
cols = cols[0:2] + cols[1:] + cols[2:3] + cols[5:6] + cols[3:5] + cols[6:-1]
df_reviews_copy.loc[:, 'UserID'] = df_reviews_copy[cols]
```

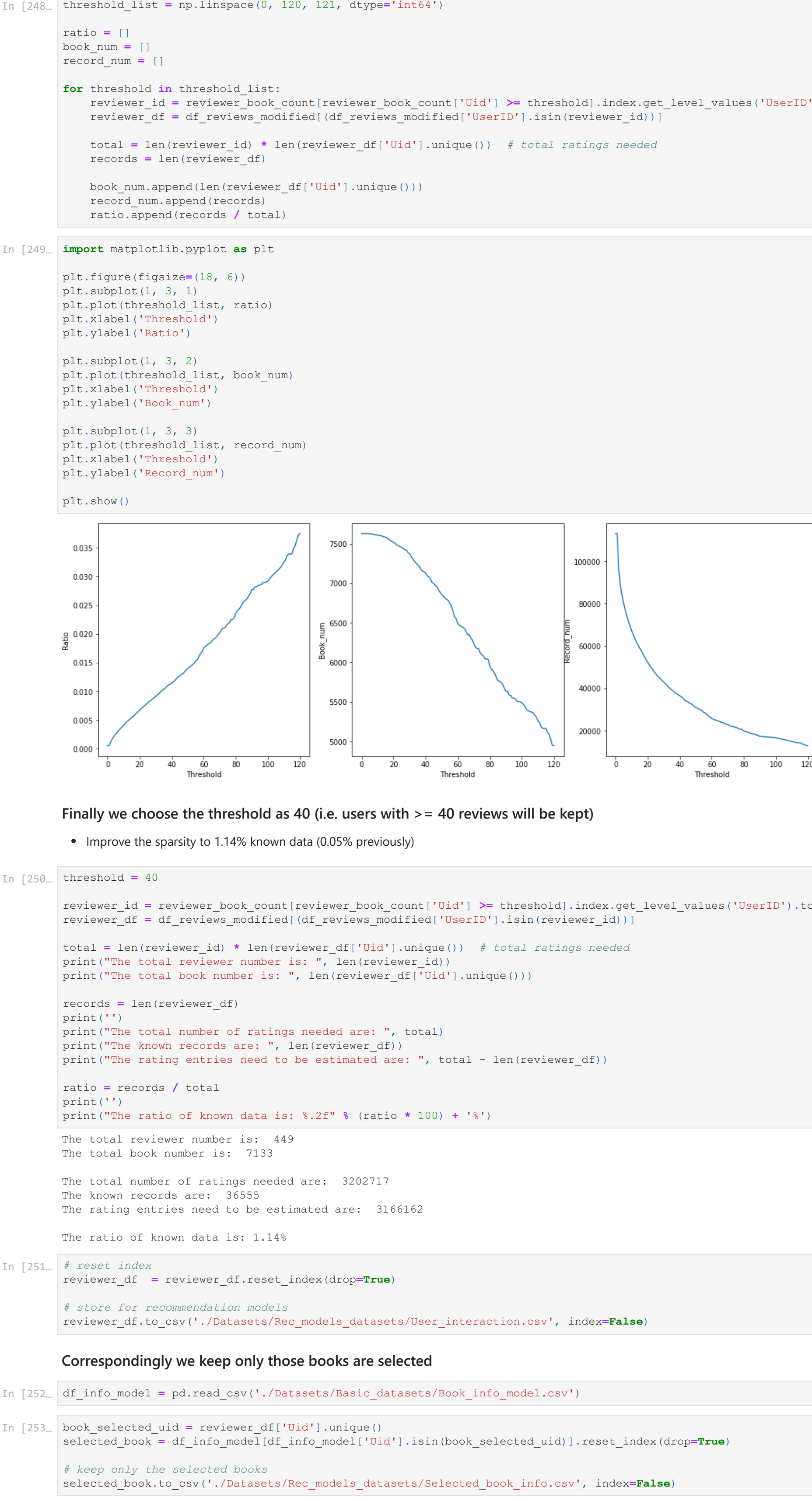
We can see that 3 different users with the same name 'Luna' has been identified and the wrong `N_Follower` has been modified

```
In [136]: df_reviews_copy[df_reviews_copy['Reviewer'] == 'Luna']

Out[136]:
```

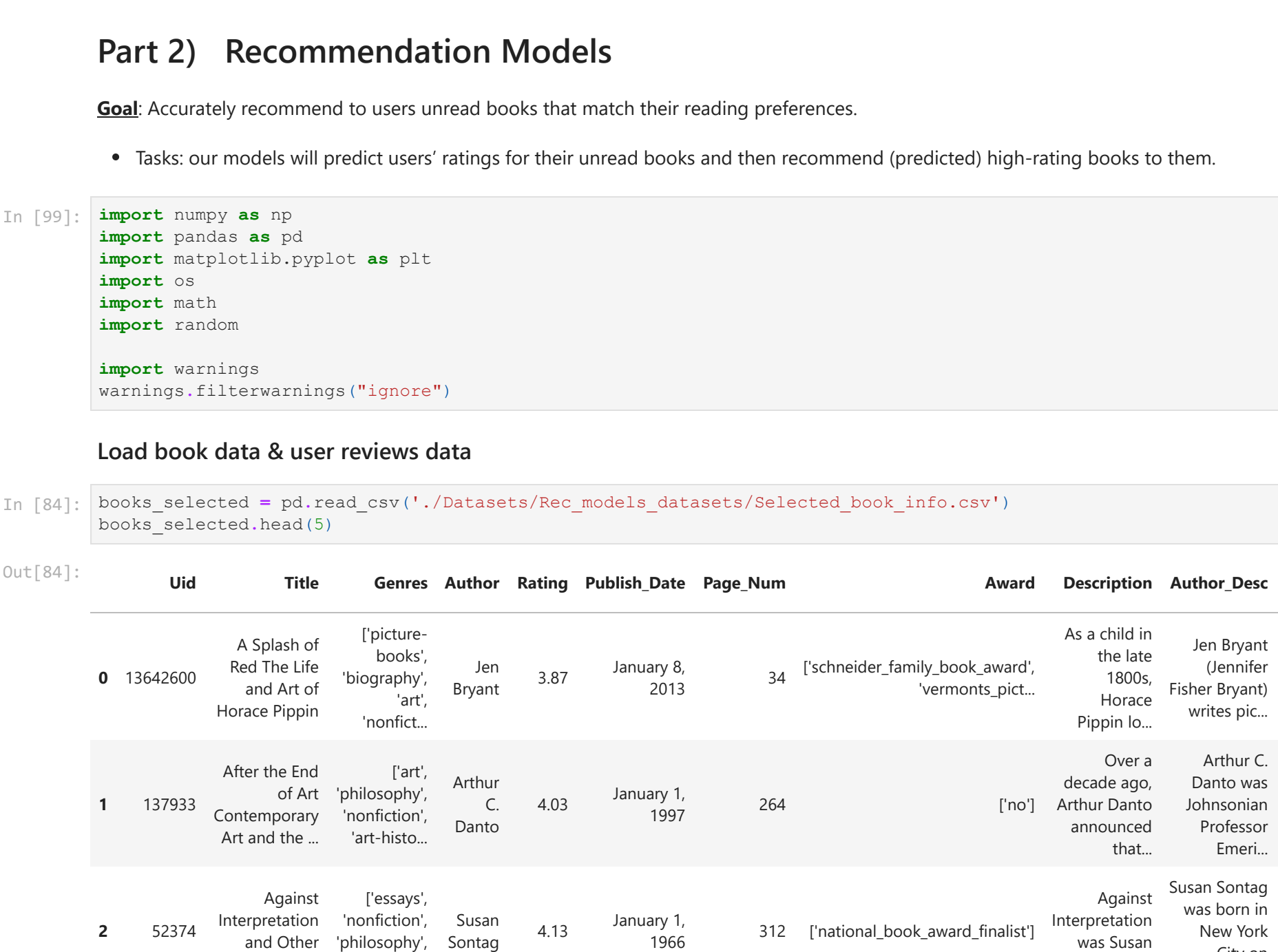
| UId   | Title  | UserID       | Reviewer | Review_Rating | N_Review | N_Follower | Review_Date | Content          | N_Likes                                          | N_Comments |    |
|-------|--------|--------------|----------|---------------|----------|------------|-------------|------------------|--------------------------------------------------|------------|----|
| 32810 | 465226 | The Redebart | 15570    | Luna          | 5        | 126        | 23          | October 28, 2022 | INCREDIBLE. Really just a work of art! You kn... | 34         | 18 |
| 36645 | 4      |              |          |               |          |            |             |                  |                                                  |            |    |



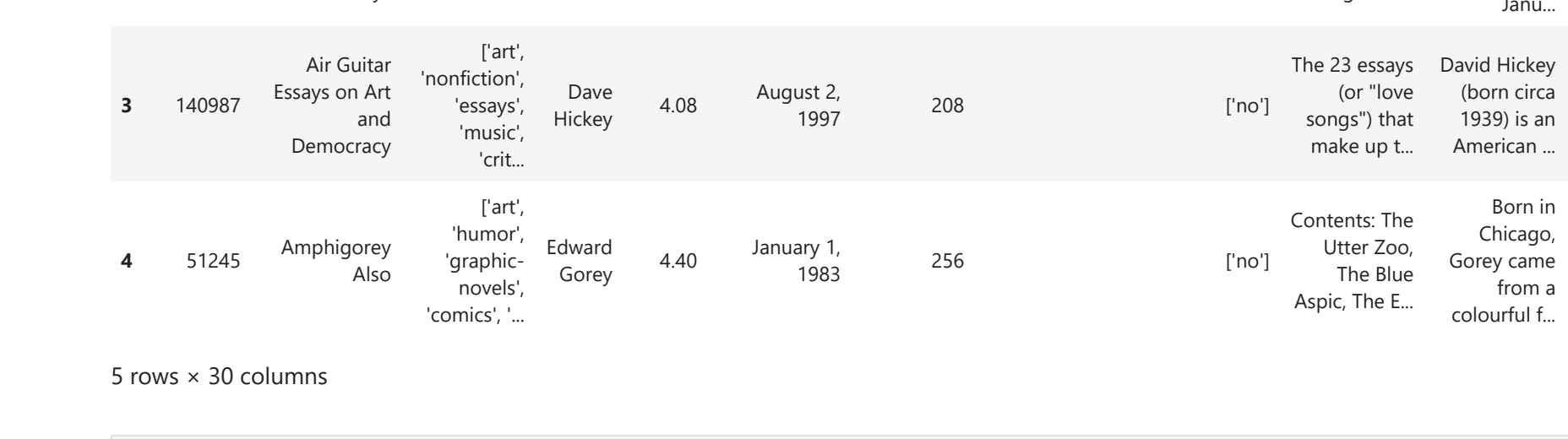


Finally we choose the threshold as 40 (i.e. users with >= 40 reviews will be kept)

- Improve the sparsity to 1.14% known data (0.05% previous)



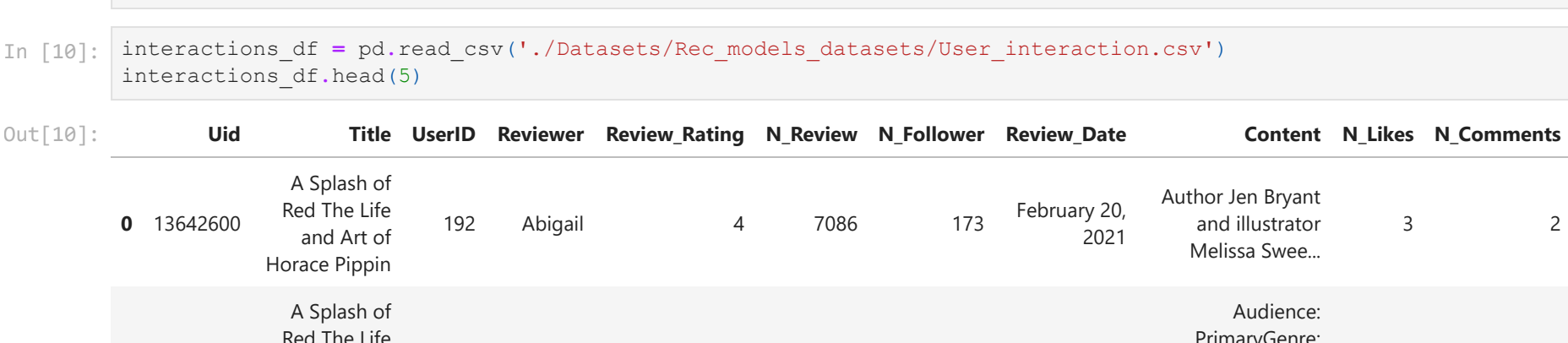
Correspondingly we keep only those books are selected



## Part 2) Recommendation Models

Goal: Accurately recommend to users unread books that match their reading preferences.

- Tasks: our models will predict users' ratings for their unread books and then recommend (predicted) high-rating books to them.



|   | UId      | Title                                             | Genres                                             | Author     | Rating | Publish Date    | Page Num | Award                                          | Description                                       | Author Desc                                |
|---|----------|---------------------------------------------------|----------------------------------------------------|------------|--------|-----------------|----------|------------------------------------------------|---------------------------------------------------|--------------------------------------------|
| 0 | 13642600 | A Splash of Red The Life and Art of Horace Pippin | ['picture-books', 'biography', 'art', 'nonfic...'] | Jen Bryant | 3.87   | January 8, 2013 | 34       | [schneider_family_book_award, vermonts_pict... | As a child in the late 1800s, Horace Pippin lo... | Jen Bryant (Jennifer Bryant) writes pic... |

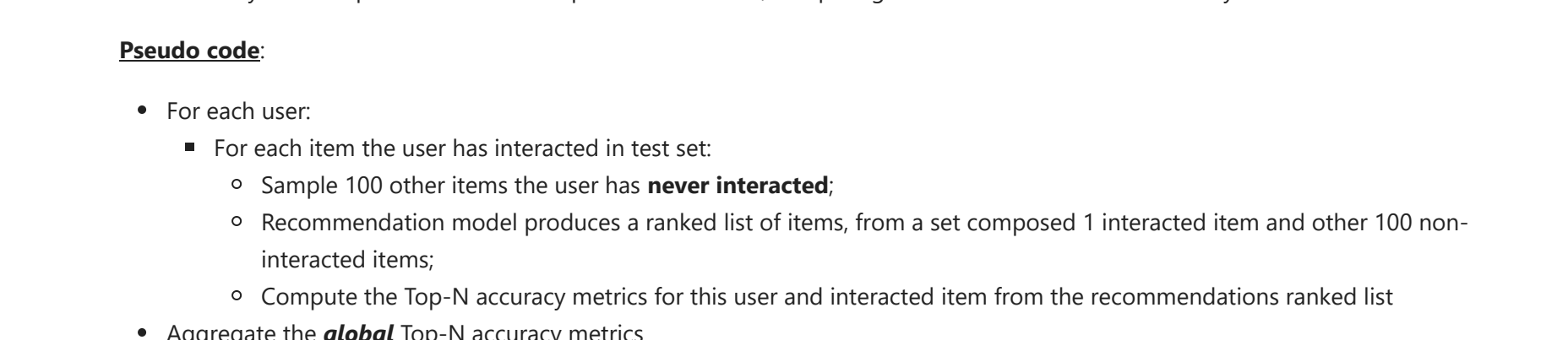
|   |        |                                                   |                                                   |                 |      |                 |     |      |                                                      |                                           |
|---|--------|---------------------------------------------------|---------------------------------------------------|-----------------|------|-----------------|-----|------|------------------------------------------------------|-------------------------------------------|
| 1 | 137993 | After the End of Art Contemporary Art and the ... | ['art', 'philosophy', 'nonfiction', 'art-histo... | Arthur C. Danto | 4.03 | January 1, 1997 | 264 | [no] | Over a decade ago, Arthur C. Danto announced that... | Arthur C. Danto was a Professor Emerit... |
|---|--------|---------------------------------------------------|---------------------------------------------------|-----------------|------|-----------------|-----|------|------------------------------------------------------|-------------------------------------------|

|   |       |                                         |                                    |              |      |                 |     |                                 |                                                  |                                                   |
|---|-------|-----------------------------------------|------------------------------------|--------------|------|-----------------|-----|---------------------------------|--------------------------------------------------|---------------------------------------------------|
| 2 | 52374 | Against Interpretation and Other Essays | ['essays', 'nonfiction', 'art...'] | Susan Sontag | 4.13 | January 1, 1966 | 312 | [national_book_award, finalist] | Against Interpretation was Susan Sontag's fir... | Susan Sontag was born in New York City on Janu... |
|---|-------|-----------------------------------------|------------------------------------|--------------|------|-----------------|-----|---------------------------------|--------------------------------------------------|---------------------------------------------------|

|   |        |                                      |                                                   |             |      |                |     |      |                                                   |                                                  |
|---|--------|--------------------------------------|---------------------------------------------------|-------------|------|----------------|-----|------|---------------------------------------------------|--------------------------------------------------|
| 3 | 140587 | An Guitar Essay on Art and Democracy | ['art', 'nonfiction', 'essays', 'music', 'crit... | Dave Hickey | 4.08 | August 2, 1997 | 208 | [no] | The 23 essays (or "love songs") that make up L... | Dave Hickey (for "love songs") that make up L... |
|---|--------|--------------------------------------|---------------------------------------------------|-------------|------|----------------|-----|------|---------------------------------------------------|--------------------------------------------------|

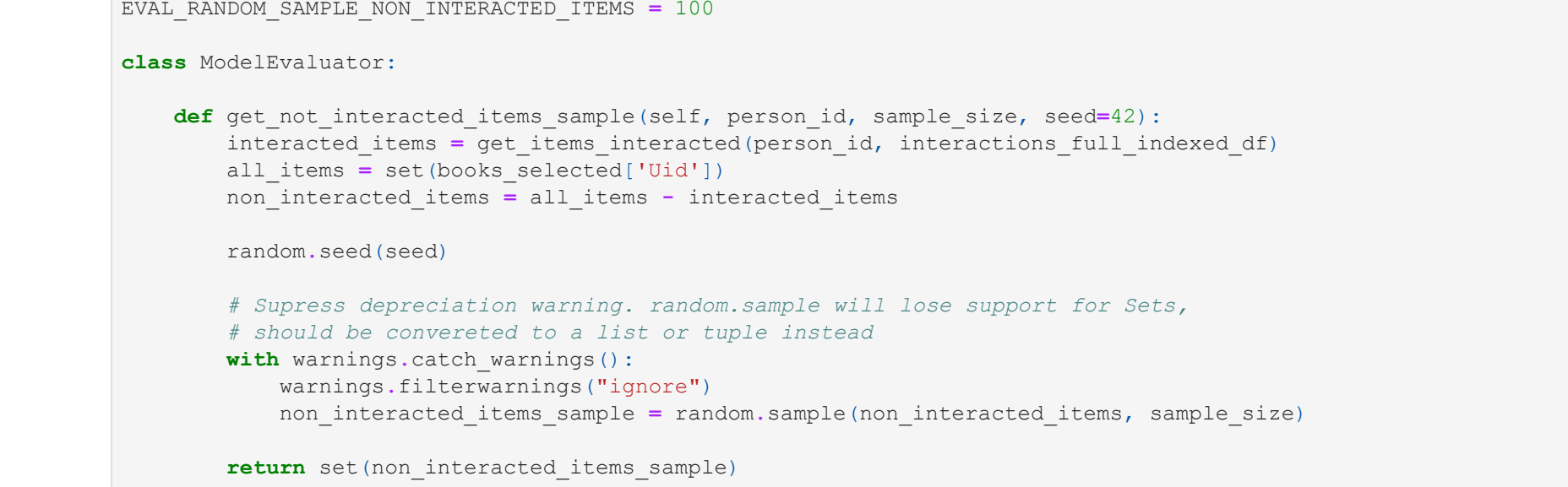
|   |       |                |                                                   |              |      |                 |     |      |                                                    |                                         |
|---|-------|----------------|---------------------------------------------------|--------------|------|-----------------|-----|------|----------------------------------------------------|-----------------------------------------|
| 4 | 51245 | Amphigory Also | ['art', 'humor', 'graphic novels', 'comics', '... | Edward Gorey | 4.40 | January 1, 1983 | 256 | [no] | Contents: The Utter Zing, The Blue Aspic, The E... | Edward Gorey came from a colourful f... |
|---|-------|----------------|---------------------------------------------------|--------------|------|-----------------|-----|------|----------------------------------------------------|-----------------------------------------|

5 rows x 30 columns



|   | UId      | Title                                             | UserId | Reviewer | Review_Rating | N_Review | N_Follower | Review_Date       | Content                                            | N_Likes | N_Comments |
|---|----------|---------------------------------------------------|--------|----------|---------------|----------|------------|-------------------|----------------------------------------------------|---------|------------|
| 0 | 13642600 | A Splash of Red The Life and Art of Horace Pippin | 192    | Abigail  | 4             | 7086     | 173        | February 20, 2021 | Author Jen Bryant and illustrator Melissa Sweet... | 3       | 2          |
| 1 | 13642600 | A Splash of Red The Life and Art of Horace Pippin | 17464  | Michelle | 4             | 301      | 38         | July 10, 2013     | Audience: Primary/Genre: Non-Fiction/Informatio... | 3       | 0          |
| 2 | 137993   | After the End of Art Contemporary Art and the ... | 17298  | Michael  | 5             | 718      | 969        | November 14, 2007 | Art itself, a great book despite my predict...     | 1       | 1          |
| 3 | 137993   | After the End of Art Contemporary Art and the ... | 13002  | Kate     | 5             | 396      | 316        | January 7, 2008   | Changed entirely how I think about art. It sta...  | 1       | 0          |
| 4 | 52374    | Against Interpretation and Other Essays           | 17298  | Michael  | 4             | 718      | 969        | March 16, 2016    | A wide-ranging debate collection of essays on a... | 67      | 2          |

## Train & Test split



## Evaluation - Top-N accuracy metrics

We use the Top-N accuracy metrics to evaluate our recommendation models' performance:

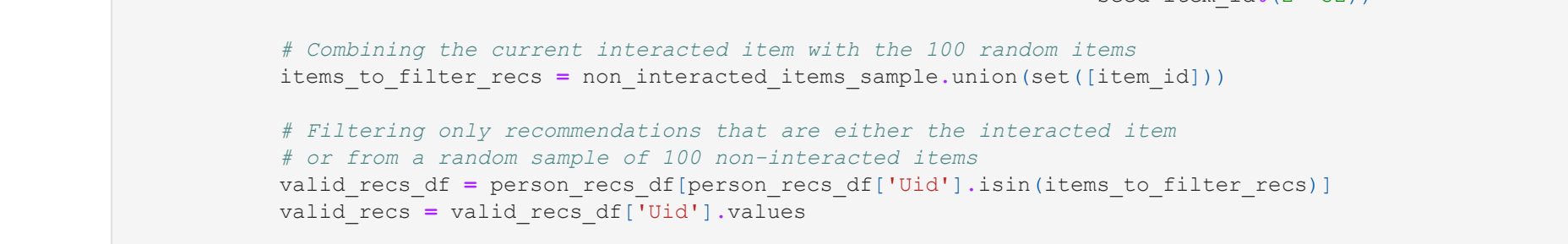
- Accuracy of the top recommendations provided to a user, comparing to the items the user has actually interacted in test set

Pseudo code:

- For each user:
  - For each item the user has interacted in test set:
    - Sample 100 other items the user has **not** interacted with
    - Recommendation model produces a ranked list of items, from a set composed 1 interacted item and other 100 non-interacted items
    - Compute the Top-N accuracy metrics for this user and interacted item from the recommendations ranked list
- Aggregate the **global** Top-N accuracy metrics

Here we use Recall@5 whether the interacted item is among the top N items (hit) in the ranked list of 101 recommendations for a user

- Metrics: Recall@5 and Recall@10
- E.g. Recall@5: For one user, if we have 100 randomly selected books, the percentage of the interacted books in the test set will be ranked among the top 5 books by the model



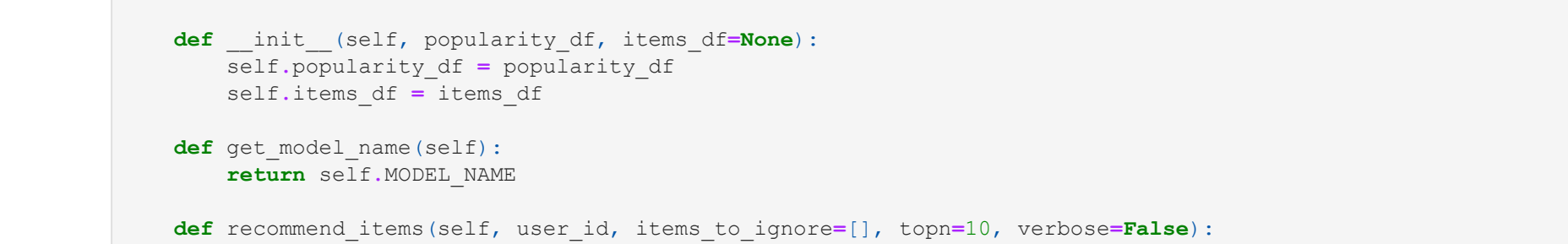
## 2.1 Popularity Model (Baseline)

This model is not actually personalized - it simply recommends to a user **the most popular items that the user has not previously consumed**.

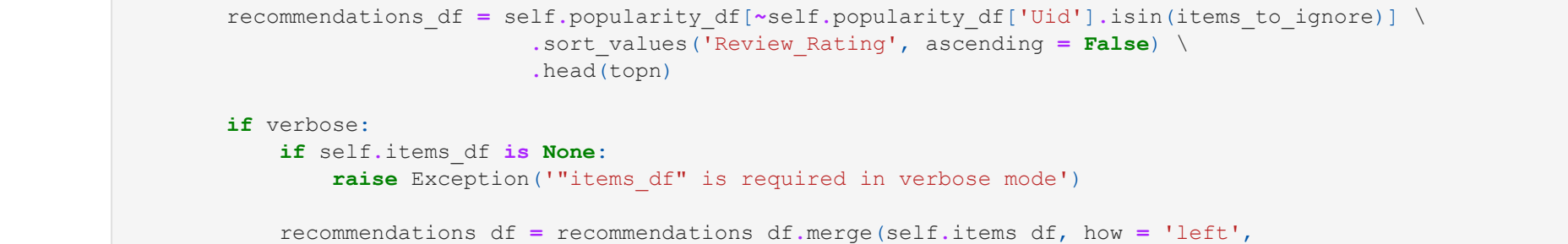
- As the popularity accounts for the 'wisdom of the crowds', it usually provides good recommendations, generally interesting for most people:
- But the main objective of a recommender system is to leverage the **long-tail items** to the users with very specific interests, which goes far beyond this simple technique. So we will use this model as a **baseline**.



|   | UId      | Review_Rating |
|---|----------|---------------|
| 0 | 36576608 | 73            |
| 1 | 464260   | 87            |
| 2 | 11588    | 71            |
| 3 | 17899948 | 71            |
| 4 | 5379586  | 69            |
| 5 | 96356    | 69            |
| 6 | 18386    | 68            |
| 7 | 7126     | 68            |
| 8 | 5907     | 67            |
| 9 | 36373463 | 67            |



## Initialize an instance for Popularity model



|     | hits@5_count | hits@10_count | interacted_count | recall@5 | recall@10 | person_id |
|-----|--------------|---------------|------------------|----------|-----------|-----------|
| 76  | 20           | 33            | 76               | 0.263158 | 0.43421   | 15602     |
| 40  | 14           | 26            | 70               | 0.200000 | 0.371429  | 17161     |
| 21  | 9            | 24            | 69               | 0.130435 | 0.347826  | 2557      |
| 7   | 10           | 20            | 63               | 0.158730 | 0.317460  | 1863      |
| 92  | 15           | 26            | 63               | 0.238095 | 0.412698  | 7740      |
| 45  | 12           | 22            | 62               | 0.193548 | 0.354839  | 26336     |
| 193 | 21           | 30            | 62               | 0.338710 | 0.483871  | 22026     |
| 57  | 14           | 16            | 61               | 0.229508 | 0.262595  | 421       |
| 267 | 21           | 31            | 60               | 0.350000 | 0.516667  | 16300     |
| 27  | 19           | 30            | 56               | 0.339286 | 0.535714  | 16765     |

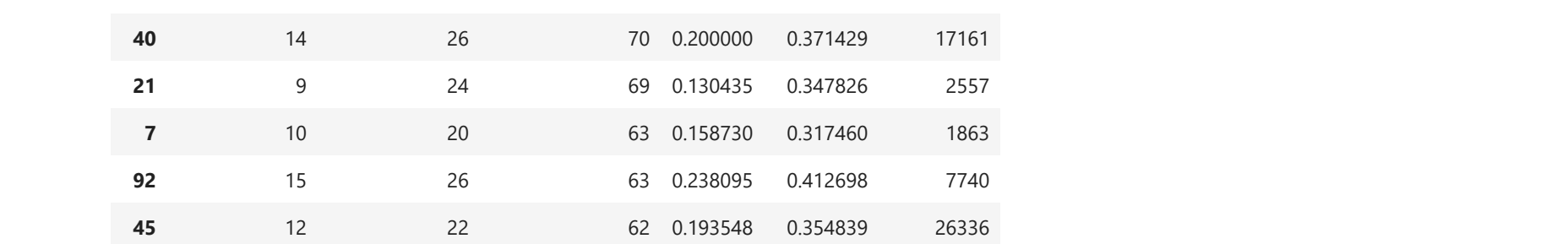
## 2.2 Content-based Filtering

Content-based filtering approaches leverage description or attributes from items the user has interacted to recommend **similar items**.

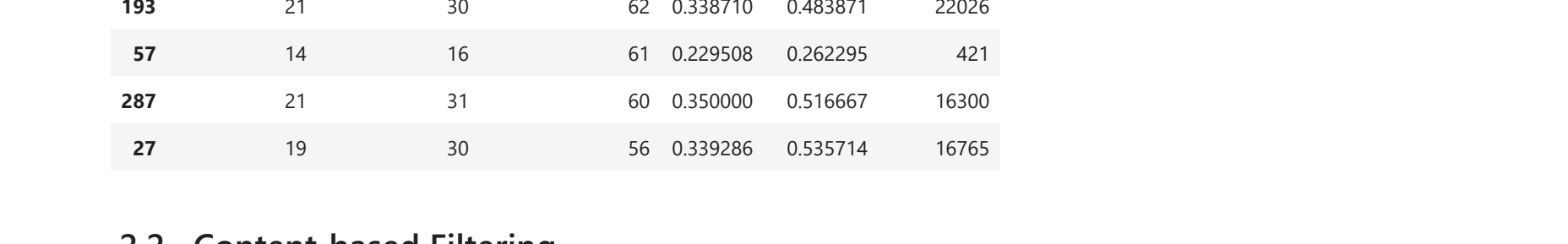
- It builds **users' profiles (tastes)** depending only on the user previous choices, making this method robust to avoid the cold-start problem

Here we are using a very popular technique in information retrieval (search engines) named **TF-IDF**.

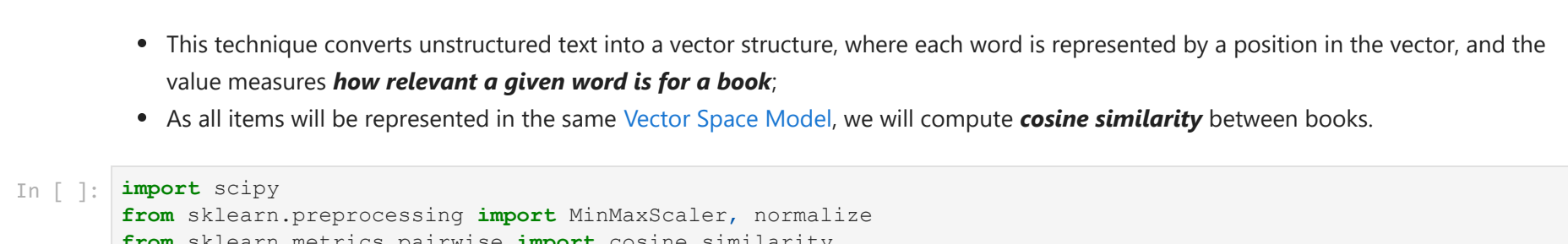
- This technique converts unstructured text into a vector structure, where each word is represented by a position in the vector, and the value measures **how relevant a given word is for a book**;
- As all items will be represented in the same **Vector Space Model**, we will compute **cosine similarity** between books.



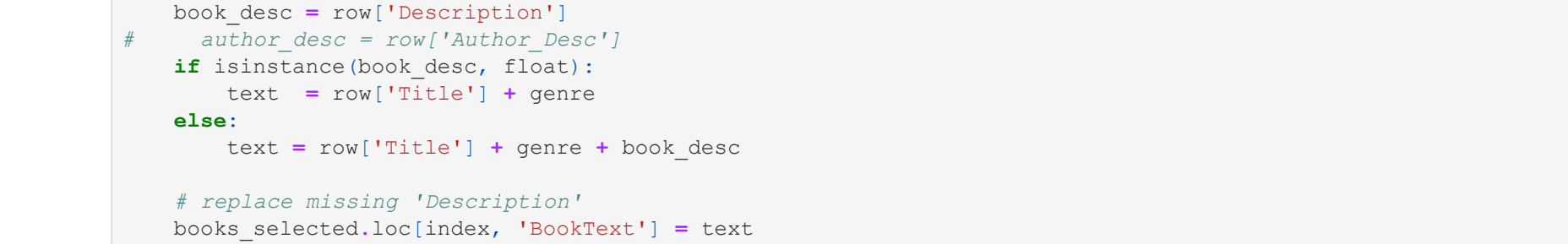
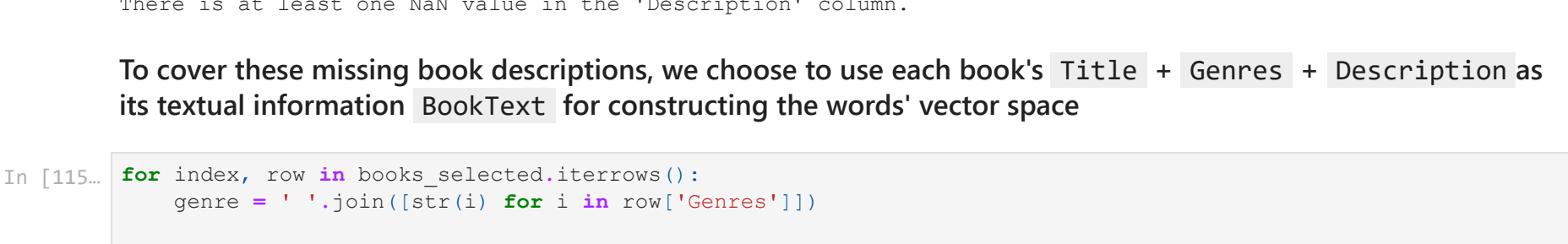
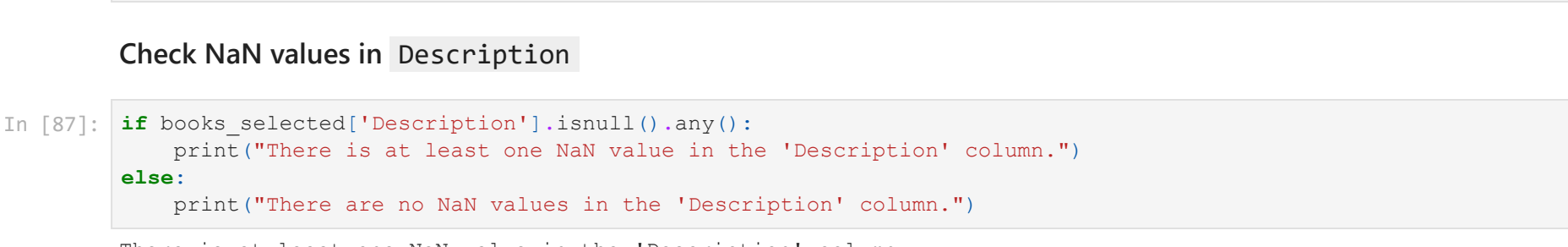
## Check NaN values in Description



To cover these missing book descriptions, we choose to use each book's **Title + Genres + Description** as its textual information **BookText** for constructing the words' vector space

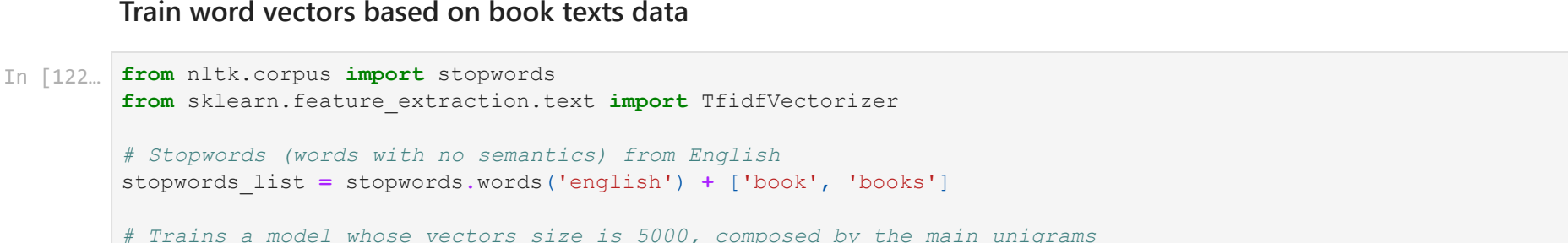


## Train word vectors based on book texts data



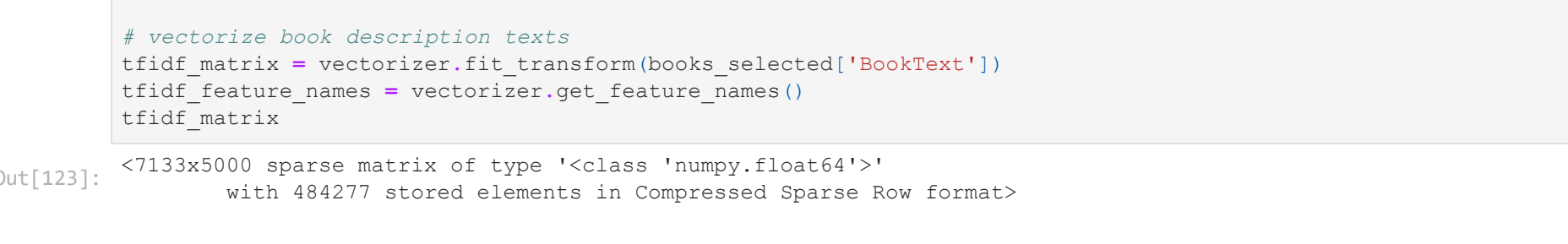
Obtain the word vector representation of a book in the same feature space as the user profiles (i.e. TF-IDF matrix)

- Take all the book files the user has interacted with and average them weighted by users' **Review\_Rating**

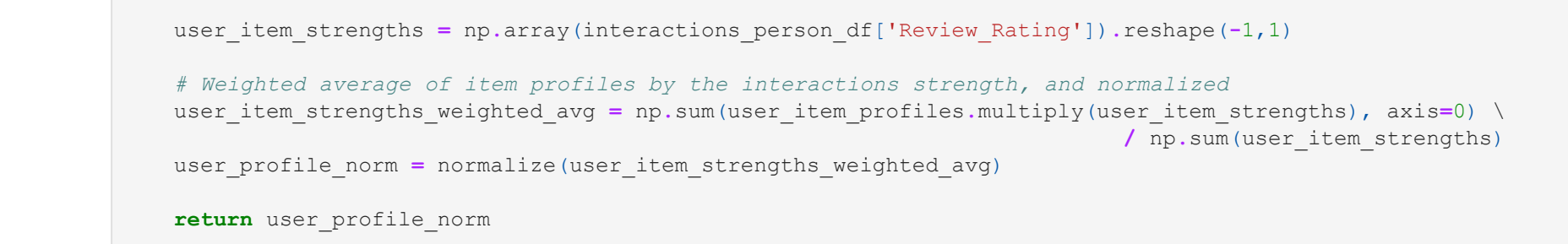
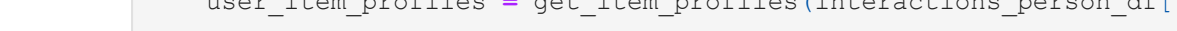
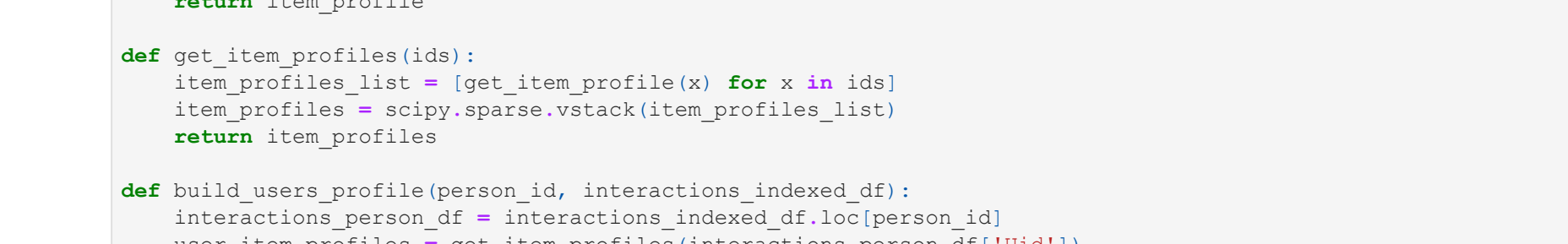


We take a look at a user profile example whose name is 'Abigail' (UserId = 192)

- The value at each position represents how relevant is a token (unigram or bigram) for the user named 'Abigail';
- It shows that she is very interested in reading books about **children**:
  - Some keywords like **bear, cat, little** appear frequently in children's books.
  - Also note that **Dr. Seuss** is a famous writer for children's literatures



|    | token            | relevance |
|----|------------------|-----------|
| 0  | children         | 0.422296  |
| 1  | dr seuss         | 0.169196  |
| 2  | seuss            | 0.169196  |
| 3  | children fiction | 0.168590  |
| 4  | animals          | 0.166772  |
| 5  | picture          | 0.157434  |
| 6  | bear             | 0.154855  |
| 7  | fiction          | 0.145393  |
| 8  | fiction classics | 0.124478  |
| 9  | dr               | 0.123242  |
| 10 | classics         | 0.112456  |
| 11 | cat              | 0.114384  |
| 12 | little           | 0.114381  |
| 13 | new              | 0.1106293 |
| 14 | hats             | 0.095753  |
| 15 | young            | 0.094332  |
| 16 | humor            | 0.094149  |
| 17 | fiction animals  | 0.090396  |
| 18 | picture fiction  | 0.089813  |
| 19 | food             | 0.088952  |



|     | hits@5_count | hits@10_count | interacted_count | recall@5 | recall@10 | person_id |
|-----|--------------|---------------|------------------|----------|-----------|-----------|
| 76  | 32           | 41            | 76               | 0.421053 | 0.539474  | 15602     |
| 40  | 21           | 34            | 70               | 0.300000 | 0.485714  | 17161     |
| 21  | 17           | 26            | 69               | 0.246377 | 0.376812  | 2557      |
| 7   | 33           | 36            | 63               | 0.523810 | 0.571429  | 1863      |
| 92  | 5            | 16            | 63               | 0.079365 | 0.253968  | 7740      |
| 45  | 19           | 26            | 62               | 0.306452 | 0.419355  | 26336     |
| 193 | 17           | 27            | 62               | 0.274194 | 0.435484  | 22026     |
| 57  | 7            | 14            | 61               | 0.114754 | 0.229508  | 421       |
| 267 | 36           | 40            | 60               | 0.600000 | 0.666667  | 16300     |
| 27  | 12           | 15            | 56               | 0.214286 | 0.267857  | 16765     |

Our content-based filtering model provides personalized recommendations with a **Recall@5 = 0.3639**, indicating that around 36% of the items that the user interacted with in the test set were included in the top-5 recommended items generated by the model from a list of 100 random items. Furthermore, the **Recall@10 = 0.4869**. This is a great improvement from baseline, partially proves the power of this technique.

## 2.3 Collaborative Filtering model

We use latent factor model (model-based) for collaborative filtering:

- Compress user-item matrix into a **low-dimensional** representation in terms of latent factors, solving sparsity problem;
- Here we use a popular latent factor model named **Singular Value Decomposition (SVD)**.

Singular Value Decomposition:

$$M_{m \times n} = U_{m \times m} * \Sigma_{m \times n} * V_{n \times n}^T \text{ (Full SVD)}$$

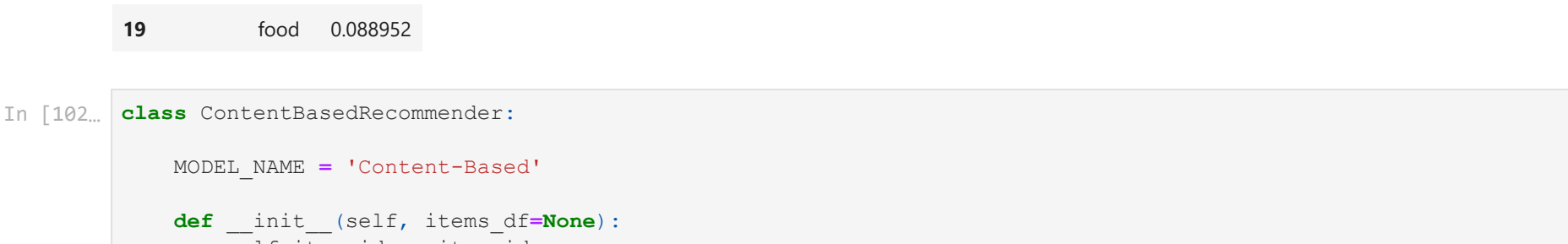
However, we can choose the **number of latent factors k** to factor the user-item matrix and approximate original matrix using the following formula:

$$M_{m \times n} \approx U_{m \times k} * \Sigma_{k \times k} * V_{k \times n}^T \text{ (Thin SVD)}$$

- The higher the number of factors, the more precise is the factorization in the original matrix reconstructions.
- But if the model is allowed to memorize too much details of the original matrix, it may not generalize well for data it was not trained on. Reducing the number of factors increases the model generalization.
- We will look over a potential space to find the **best number of factors**  $k_{optimal}$



## Matrix Factorization









```

In [169]: from collections import OrderedDict
import numpy as np
import spacy
from spacy.lang.en.stop_words import STOP_WORDS
import string
import string
from string import punctuation

nlp = spacy.load('en_core_web_sm')

class TextRankKeyword():
 """Extract keywords from text"""

 def __init__(self):
 self.d = 0.85 # damping coefficient, usually is .85
 self.min_diff = 1e-5 # convergence threshold
 self.steps = 10 # iteration steps
 self.node_weight = None # save keywords and its weight

 def set_stopwords(self, stopwords):
 """Get stop words"""
 for word in STOP_WORDS.union(set(stopwords)):
 lexeme = nlp.vocab[word]
 lexeme.is_stop = True

 def sentence_segment(self, doc, candidate_pos, lower):
 """Store those words only in candidate_pos"""
 sentences = []
 for sent in doc.sents:
 selected_words = []
 for token in sent:
 # Store words only with candidate POS tag
 if token.pos_ in candidate_pos and token.is_stop is False:
 if lower is True:
 selected_words.append(token.text.lower())
 else:
 selected_words.append(token.text)
 sentences.append(selected_words)
 return sentences

 def get_vocab(self, sentences):
 """Get all tokens"""
 vocab = OrderedDict()
 i = 0
 for sentence in sentences:
 for word in sentence:
 if word in ('.', ',', '!', '?', '-', '/', '\n'):
 continue
 if word not in vocab:
 vocab[word] = i
 i += 1
 return vocab

 def get_token_pairs(self, window_size, sentences):
 """Build token pairs from windows in sentences"""
 token_pairs = list()
 for sentence in sentences:
 for i, word in enumerate(sentence):
 if word in ('.', ',', '!', '?', '-', '/', '\n'):
 continue
 for j in range(i+1, i+window_size):
 if j >= len(sentence):
 break
 pair = (word, sentence[j])
 if pair not in token_pairs:
 token_pairs.append(pair)
 return token_pairs

 def symmetrize(self, a):
 return a + a.T # np.diag(a.diagonal())

 def get_matrix(self, vocab, token_pairs):
 """Get normalized matrix"""
 # Build matrix
 vocab_size = len(vocab)
 g = np.zeros((vocab_size, vocab_size), dtype='float')
 for word1, word2 in token_pairs:
 i, j = vocab[word1], vocab[word2]
 g[i][j] = 1
 # Get Symmetric matrix
 g = self.symmetrize(g)
 # Normalize matrix by column
 norm = np.sum(g, axis=0)
 g_norm = np.divide(g, norm, where=norm!=0) # this is ignore the 0 element in norm
 return g_norm

 def get_keywords(self, number=10):
 """Print top number keywords"""
 keywords_list = []
 node_weight = OrderedDict(sorted(self.node_weight.items(), key=lambda t: t[1], reverse=True))
 for i, (key, value) in enumerate(node_weight.items()):
 if key not in list(punctuation):
 keywords_list.append(key)
 if i > number:
 break
 return keywords_list

 def analyze(self, text,
 candidate_pos=['NOUN', 'ADJ', 'VERB'],
 window_size=4, lower=False, stopwords=list()):
 """Main function to analyze text"""
 self.set_stopwords(stopwords)
 # Pare text by spaCy
 doc = nlp(text)
 # Filter sentences
 sentences = self.sentence_segment(doc, candidate_pos, lower) # list of list of words
 # Build vocabulary
 vocab = self.get_vocab(sentences)
 # Get token_pairs from windows
 token_pairs = self.get_token_pairs(window_size, sentences)
 # Get normalized matrix
 g = self.get_matrix(vocab, token_pairs)
 # Initialization for weight (pagerank value)
 pr = np.array([1] * len(vocab))
 # Iteration
 previous_pr = 0
 for epoch in range(self.steps):
 pr = (1-self.d) + self.d * np.dot(g, pr)
 if abs(previous_pr - sum(pr)) < self.min_diff:
 break
 else:
 previous_pr = sum(pr)
 # Get weight for each node
 node_weight = dict()
 for word, index in vocab.items():
 node_weight[word] = pr[index]
 self.node_weight = node_weight

We try with example book (Uid = 449128): Alla Prima: Everything I Know about Painting

- Window size k = 4, only consider NOUN & ADJ words
- We can see it generally catches important words like: artist, painting, canvas, color, etc.

In [170]: df_comments.iloc[5, 0]
449128

Out[170]:

In [171]: text = df_comments.iloc[5, 1]
tr4w = TextRankKeyword()

We only want adjective & Nouns for tags
tr4w.analyze(text, candidate_pos = ['ADJ', 'NOUN'], window_size=4, lower=False)
tr4w.get_keywords(10)

Out[171]:
['artist',
 'advice',
 'color',
 'schmid',
 'painters',
 'confidence',
 'oil',
 'painting',
 'painter',
 'canvas',
 'artists',
 'composition']

We can also get each keyword's weights

In [173]: [key: tr4w.node_weight[key] for key in tr4w.get_keywords(10)]

Out[173]:
[('pippin': 9.081010937634829),
 ('horace': 7.301031593834669),
 ('artist': 3.3053121859984687),
 ('red': 2.849867898162826),
 ('pippin': 2.821484620122126),
 ('illustrations': 2.80323385238113),
 ('artists': 2.802856402697573),
 ('works': 2.478942374657984),
 ('painter': 2.3143430205118087),
 ('composition': 2.2843430205118087)]

```

We can also get each keyword's weights

- ```
[[key: tr4w.node_weight[key]] for key in tr4w.get_keywords(10)]
```
- ```
{'pippin': 9.081910937634829},
{'horace': 7.501031593836607},
{'artist': 3.9051219599946987},
{'red': 2.845867898168256},
{'pippins': 2.8214848201212126},
{'...': ...}
```

```

['illustrations': 2.262229982361193,
 'artists': 2.8028856402699793,
 'artwork': 2.478949233735457984,
 'painter': 2.2143430305118087,
 'paintings': 2.27211299262405882,
 'pictures': 2.262719780963017,
 'notes': 2.0621940091953608]

```

Next we used *TexRank* to extract 5 keywords of each book and classify each review to that tags if it contains the keyword

- The following process is time-consuming, just use our stored result datasets
  - book\_tags.csv** 5 keywords tags for each book
  - review\_tags.csv** keywords this review contains

```

In []: # count how many reviews each book has
reviews_count = reviews.groupby(['Uid'], sort=False)['Uid'].count()

add column 'Tags' if this review contains the keyword
reviews['Tags'] = []
df_comments['Tags'] = []

num = 0 # the number of reviews we have classified
for i in tqdm(range(df_comments.shape[0])):
 print("The %d-th book" % i)
 text = df_comments.iloc[i, 1] # loop over each book

 # get 5 keywords
 tr4w = TextRankKeyword()
 tr4w.analyze(text, candidate_pos = ['ADJ', 'NOUN'], window_size=4, lower=False)
 keyword_list = tr4w.get_keywords(5)

 # for each book, store their keywords from reviews
 df_comments['Tags'][i] = keyword_list

 # for each review, assign their tags if containing the keyword
 for in_range(reviews_count.values[i]):
 tag_list = []
 for key in keyword_list:
 if key in reviews.iloc[num,1].split():
 tag_list.append(key)

 if not tag_list:
 num += 1
 continue
 print(tag_list)
 reviews['Tags'][num] = tag_list

```

```

 if not tag_list:
 num += 1
 continue
 print(tag_list)
 reviews['Tags'][num] = tag_list

```

- ```
num = num + 1
```
- We look at the tag extracted results
- ```
each book's keywords (tags)
df_comments.head(8)
```
- | UId | Content |
|-----|---------|
|-----|---------|

```

0 13642600 beautifully illustrated horace pippin finish... [pippin, horace, artist, red, pippins, illustr...
1 470185 strangest moments suddenly appear everywhere k... [animals, essays, animal, photography, berger,...
2 137933 philosophers tend worst theorist artists tend... [artists, philosophical, artist, criticism, co...
3 52374 wideranging debut essays film stimulating rel... [sontag, interpretation, essays, criticism, in...
4 140987 introduced draw hickey painting classes profe... [basketball, essays, hickey, criticism, hickey...
5 449128 painter master artists richard schmid artists... [artist, advice, color, schmid, painters, conf...
6 51245 included utter zoothe blue aspicthe epiplecti... [gorey, goreys, legacy, blue, loathsome, uniter...
7 47559 convey joy edward goreys adore unstung harp... [gorey, hapless, times, drawings, gashlycrumb...

```

---

```

In [108]: # each review's tags (if containing the book's keywords)
 reviews.head(8)

```

---

```

Out[108]:
```

|   | Uid      | Content                                           | Tags                                             |
|---|----------|---------------------------------------------------|--------------------------------------------------|
| 0 | 13642600 |                                                   | []                                               |
| 1 | 13642600 | beautifully illustrated horace pippin finish g... | [pippin, horace, artist, ...]                    |
| 2 | 13642600 | featured grandma reads session splash red hor...  | [pippin, horace, red, artists]                   |
| 3 | 13642600 | cute nonfiction artist horace pippin grad clas... | [pippin, horace, artist]                         |
| 4 | 13642600 | splash red horace pippin wellresearched pictur... | [pippin, horace, red, illustrations]             |
| 5 | 13642600 | jen bryant illustrator melissa collaborated pi... | [pippin, horace, artist, red, artists]           |
| 6 | 13642600 | audience primarygenre hisher talents splash re... | [pippin, horace, artist, red, illustrations]     |
| 7 | 13642600 | aloud bryant obligations frequently artist rep... | [pippin, horace, artist, pippins, illustrations] |

---

```

In [197]: # store for use
 df_comments.to_csv('output/book_tags.csv')
 reviews.to_csv('output/review_tags.csv')

```

---

### 3.4 Other Trials: Key Phrases

We have also considered using **key phrases** instead of keywords as our tags and we have tried many methods like **Yake1**, **KeyBERT**, **Rake** etc. to extract keyphrases, but the phrases we got are not very meaningful as a tag.

We guess for these methods to get a better result, they might need more effort for tuning model parameters and improve the quality of original textual reviews, both could be time-consuming and no promise for final result. Therefore, we conclude these methods are not very suitable and finally reject them.

---

```

In [193]: # previous tag result from TextRank
 df_comments.iloc[3, 2]

```

---

```

Out[193]: 'sontag',
 'interpretation',
 'essays',
 'criticism',
 'intellectual',
 'essay',
 'content']

```

---

### 3.4.1 YAKE!

YAKE! is a light-weight unsupervised automatic keyword extraction method which rests on text **statistical features** extracted from single documents to select the most important keywords of a text.

---

```

In [185]: import yake

```

YAKE! is a light-weight unsupervised algorithm for selecting keywords from documents to select the most important keywords.

```
In [18]: import yake

Take one example
text = df.comments.iloc[3,1]
kw_extractor = yake.KeywordExtractor(topn=10, n = 2, stopwords=None)
keywords = kw_extractor.extract_keywords(text)

for kw, v in keywords:
 print("Keyphrase: ", kw, " : score:", v)

Keyphrase: notes camp : score 0.0001778716562028035
Keyphrase: susan sontag : score 0.0005297176952421673
Keyphrase: essays susan : score 0.000722464313954816
Keyphrase: interpretation essays : score 0.0007763338108946728
Keyphrase: sontag argues : score 0.0008288173643892819
Keyphrase: films robert : score 0.000846658595228895
Keyphrase: saetres saint : score 0.0008950281626283484
Keyphrase: saint gene : score 0.0009264273338927114
Keyphrase: imagination disaster : score 0.001005120404112312
Keyphrase: interpretation sontag : score 0.00112568756802986

According to the result, we could see there are some repetitions and overlaps across the phrases, which makes them hard to convey more
diverse meanings and cause information redundancy, which is bad for tags selection.
```

### 3.4.2 TextRank (Summa)

The core method is the same as TextRank we used before. But this package is not that intelligent since there are so many repetition and it did not catch important information of book reviews.

```
In [186]: from summa import keywords

text = df.comments.iloc[3,1]
TR_keywords = keywords.KeywordRanker(text, score=True)
print(TR_keywords.get(0)[0])

('!sontag', 0.407184696812149), ('sontag', 0.47001946696812149), ('interpretation', 0.2987048042299369), ('critic
3', 0.2654791090129931), ('criticized', 0.2654791090129931), ('essays film stimulating', 0.1835148264138276
3), ('interpretations cultured', 0.16746065948308786), ('interpret impoverish', 0.16343202148891145), ('interpret
1 taken down', 0.14861324210370381), ('essay notes camp remain', 0.14395467616787973)]
```

### 3.4.3 KeyBERT

KeyBERT is a simple, easy-to-use keyword extraction algorithm that takes advantage of SBERT embeddings to generate keywords and key

keyphrase is a `tfidf` embedding. The `tfidf` embedding algorithm (takes advantage in `tfidf` embeddings to generate keywords and phrases from a document that are more similar to the document).

- First, document embedding (a representation) is generated using the sentences-BERT model;
- Next, the embeddings of words are extracted for N-gram phrases. Computed cosine similarity of each keyphrase to the document;
- The most similar words can then be identified as the words that best describe the entire document and are considered as keywords.

Since BERT model need GPU, we put the result in the comments below.

```

In []: from keybert import KeyBERT

text = df_comments.iloc[3,1]
kw_model = KeyBERT(model='all-mpnet-base-v2')
keywords = kw_model.extract_keywords(text,
 keyphrase_ngram_range=(1, 3),
 stop_words='english',
 highlight=False,
 top_n=10)

keywords_list= list(dict(keywords).keys())
print(keywords_list)

keywords_list \ ['pippin art', 'pippin paintings', 'artist horace', 'pippin artist', 'horace pippin', 'pippin impressionistic', 'painter horace',
'evocative pippin', 'art horace', 'watercolor paintings']

```

The phrases are not very meaningful as tags because we can see too many common parts between these phrases (e.g. pippin), which might not be able to comprehensively reflect a book's attributes as tags.

### 3.4.4 Rake

Rake is short for **Rapid Automatic Keyword Extraction** - a method of extracting keywords from individual documents. It can also be applied to new fields very easily and is very effective in dealing with multiple types of documents, especially text that requires specific grammatical conventions.

Rake identifies key phrases in a text by analyzing the occurrence of a word and its compatibility with other words in the text (co-occurrence).

```

In [195]: from multi_rake import Rake

```

occurrence).

```
[195]: from multi_rake import Rake

text = df_comments.iloc[3,1]
rake = Rake()
keywords = rake.apply(text)
print(keywords[10])

[('quarter pretty ephemeral', 9.0), ('packaged sontag', 4.0), ('elevated', 1.0)]
```

We can see that most phrases are not very summative, the result is far away from satisfying.

### 3.5 Extract Summary for Reviews on Each Book

After we generate a list of recommended books, we want to help our book readers get a snapshot of what are the major opinions of other people who have reviewed the book. Thus, we aim to provide a summary of reviews on each book to help readers acquire the main ideas efficiently. We utilize **extractive summarization** method to extract top sentences from all reviews of a certain book and consider the top sentences as the summary. We explored the following two methods to extract the summary:

- **Method 1:** Baseline - LexRank + Sentence-BERT embedding
- **Method 2:** Improvement based on Method 1 to achieve three objectives
  - Maximize Centrality
  - Minimize Redundancy
  - Balance Sentiment

#### 3.5.1 Method 1 - Baseline Model

- We firstly use **Sentence-BERT** to graph sentence embeddings of each sentence
- Next, we use **LexRank** to get the graph-based "centrality" of each sentence. The higher the centrality score, the more prominent the sentence is to the whole review.
- Since **cosine similarity** allows us to approximate how similar two sentence vectors are by simply using the cosine of the angle

- Next, we use **LexRank** to get the graph

Sentence is to the **whole review**.

- Since **cosine similarity** allows us to approximate how similar two sentence vectors are by simply using the cosine of the angle between the two vectors to quantify how similar two sentences are, we calculate the cosine similarity of sentence embeddings from **Sentence-BERT**.
- If one sentence is similar to many other sentences in the review, we can claim that it is more central to the review or the "center" of the graph and thus it is more important.

In summary, for baseline model, we calculate **cosine similarity** of sentence embeddings from **Sentence-BERT** and extract **5 sentences** with the highest centrality scores as the summary.

Note: Our original code files need to run with GPU, thus we only put some code snippets and the results here as demonstration for the report.

```

In []: from sentence_transformers import SentenceTransformer, util
 from tqdm import tqdm
 from sklearn.metrics.pairwise import cosine_similarity
 from nltk.tokenize import word_tokenize
 from textblob import TextBlob

 model = SentenceTransformer('sentence-transformers/paraphrase-mpnet-base-v2')

```

### Build a pipeline function for generating summary of book reviews

```

In []: def extract_summary(UID, topk=5):
 # select all review sentences of a typical book with selected UID

```

```
def extract_summary(lid, topk=5):
```

```

select all review sentences of a typical book with selected Oid
sentences = sent_tokenize(id_comments[id_comments['Oid'] == Oid].Content)

We calculate the sentiment score using TextBlob, this measurement will be used as a evaluation metric
sentiments = list(map(lambda text: TextBlob(text).sentiment.polarity, sentences))
doc_sentiment = np.mean(sentiments)

Encode each sentence with Sentence-BERT
sentence_embeddings = model.encode(sentences)

Calculate the cosine similarity scores among sentences
cos_scores = util.cos_sim(sentence_embeddings, sentence_embeddings).numpy()

Obtain the centrality scores with Louvain
centrality_scores = degree_centrality_scores(cos_scores, threshold=None)

We obtain the indexes of the most central sentences
most_central_sentence_indexes = np.argsort(-centrality_scores)

summary_indexes = list(most_central_sentence_indexes[:topk])

We use redundancy score as another evaluation metric
def compute_redundancy(idx, summary_indexes):
 if not summary_indexes: return 0
 return max([cos_scores[idx][sent] for senti in summary_indexes])

Store the outputs
out_summary = ""
out_sentiment = ""
out_centralities = ""
summary_sentiments = []
summary_centralities = []
summary_redundancies = []

for idx in summary_indexes:
 summary_sentiments.append(sentiments[idx])
 summary_centralities.append(centrality_scores[idx])

compute redundancy of the specific sentence to all other sentences in the summary
summary_indexes.remove(idx)
redundancy = compute_redundancy(idx, summary_indexes)
summary_indexes.append(idx)
summary_redundancies.append(redundancy)

```

```
compute redundancy of t
summary_idxes.remove(idx)
redundancy = compute_redundancy(t, summary_idxes)
```

```
summary_indices.append(idx)
summary_redundancies.append(redundancy)

out_summary += sentences[idx] + "[SEP]"
out_sentiment += str(sentiments[idx]) + "[SEP]"
out_centralities += str(centralities_scores[idx]) + "[SEP]"

mean_sentiment = np.mean(summary_sentiments)
mean_centrality = np.mean(summary_centralities)
mean_redundancy = np.mean(summary_redundancies)

print("(*=mean_sentiment), (mean_centrality), (mean_redundancy), (doc_sentiment)")

return out_summary, out_sentiment, out_centralities, mean_sentiment, mean_centrality, mean_redundancy, doc
```

## Result demonstration

In [28]:

```
from IPython.display import Image, display
display(Image('!.\\pictures\\Baseline_output.png', width=1200, height=1200))
```

Summary:

A picture book kido about self-taught painter Horace Pippin. Score: 1.465276808382623

A splash of red: The Life and of Horace Pippin is a well-researched picturobook biography of an esteemed self-taught African American painter. Since I knew nothing about Horace Pippin or his art before, I appreciated Jan Bryant's informative narrative. Score: 1.63208626404992

Although the text is relatively brief, it contains quite a bit of information on Horace Pippin's life and work. From his youth, to his service in WWI, to his later success as a painter. Score: 1.47086015781832

A splash of red: The Life and Art of Horace Pippin by Jan Bryant and Illustrated by Melissa Sweet. NOTES: The colors are amazing. Score: 1.39879223564607

We learn Horace Pippin was a great artist! Score: 1.36394968954632

Grateful for the authors and illustrators who bring forward people, places, and experiences we might have missed in our patchwork education. Cute notification picture book about the work of artist Horace Pippin. Score: 1.36577702955398

It is a biography of Horace Pippin, an artist that I was unfamiliar with, but also a bit of a 'coding of age' book too as Horace tries different occupations before and after he goes on before coming to the fact that he is, indeed, an artist. Score: 1.36093819713189

Review: Laraine enjoyed this beautifully illustrated book of the life of Horace Pippin! Score: 1.34383869260029

Eventually, his work was discovered and grown up with such figures as painter M. C. W. Hecht, and he went on to become a well-known painter. I am as glad that I picked up a splash of red: The Life and Art of Horace Pippin, as I had not previously heard of this artist, but will now seek out more of his work.

It is a biography of Horace Pippin, an artist  
occupations before and after the war before a  
review later. We enjoyed this beautifully illus

Eventually, his work was discovered and purchased by such figures as painter M.C. Wyeth, and he went on to become a well-known painter. "I am so very glad that I picked up a Splash of Red: The Life and Art of Horace Pippin, as I will not previously heard of this artist, but I will now seek out more of his work. Score: 1.3356312358645878

A gorgeous picture book about the life of Horace Pippin. Score: 1.393126407692999

However, after running the baseline model, we noticed some issues.

- First, there may be **different sentiment distribution** of extracted summary and original review texts of one book
  - NOTE: Sentiment scores are calculated using TextBlob Polarity Score (-1: Negative to 1: Positive)
  - For example, we may have mean sentiment of summary as 0.1565 and mean sentiment of original texts as 0.1217
- Second, there may be some **redundant information** among sentences in extracted summary
  - For example, for two sentences in the summary, they convey similar meanings
    - Eg. Good book with clear structure! & Great book with organized structure!
  - We want to try avoiding such redundancy and providing more diverse information to end users with the summary

To deal with the issues, we propose a more **balanced** summarization algorithm in Method 2.

## 3.5.2 Method 2 - Improvement - Balanced Summarization

We use a greedy algorithm to extract a more balanced summary with **three objectives**

- Maximize centrality score** which represents the importance of the sentence
- Minimize the difference between the sentence centrality and the overall centroid of the original review texts**

We use a greedy algorithm to extract a m

- **Maximize centrality score** which represents the importance of the sentence
- **Minimize** the difference between the summary sentence and the overall sentiment of the original review texts
- **Minimize** the redundant information between extracted sentences

The learning objective of the model can be written as:

$$O(s, S, A) = Centrality(s) + CosineSimilarity(s, S) - SentimentDifference(S \cup s, A)$$

where

- $s$  represent the current review sentence,
- $S$  is the target extracted summary, and
- $A$  represents all the review texts of that specific book

Essentially, we want to extract a summary with high centrality, low redundancy, and a balanced sentiment.

The following is the code snippets of our improved algorithm.

```
In []: def extract_summary_balanced(Uid, topk=5):
 """
 :param Uid:
 :param topk:
 :return: out_summary: summary string separated by [SEP]
```

```

:param Uid:
:param topic:
:return: out_summary: summary string separated by [SEP]
 out_sentiment: sentiments string separated by [SEP]
 out_centralities: centrality score string separated by [SEP]
 mean_sentiment,
 mean_centrality,
 max_redundancy
 doc_sentiment: mean sentiment of original text score

select sample review sentences
sentences = sent_tokenize(d_comments['comments']['Uid'] == Uid, Content)
calculate the sentiment score using TextBlob, this measurement will be used as a evaluation metric
sentiments = list(map(lambda text: TextBlob(text).sentiment.polarity, sentences))
doc_sentiment = np.mean(sentiments)

Encode each sentence with Sentence-BERT
sentence_embeddings = model.encode(sentences)
Calculate the cosine similarity scores among sentences
cos_scores = util.cos_sim(sentence_embeddings, sentence_embeddings).numpy()
Obtain the centrality scores with LexRank
centrality_scores = degree_centrality_scores(cos_scores, threshold=None)

We use redundancy score as another evaluation metric
def compute_redundancy(idx, summary_idxes):
 if not summary_idxes: return 0
 return max([cos_scores[idx][sentid] for sentid in summary_idxes])

We compute the sentiment score difference between the extracted summary and the original review text

```

```
def _compute_redundancy(idx,
 if not summary_idxes: ret
```

```

 return max([cos_scores[idx][sent1] for sent1 in summary_idxes])

We compute the sentiment score difference between the extracted summary and the original review text
def sentiment_difference(idx, summary_idxes):
 return abs(doc_sentiment - np.mean([sentiments[idx] + [sentiments[sent1] for sent1 in summary_idxes]))

summary_idxes = []
best_idx = None

while len(summary_idxes) < topk and best_idx != -1:
 best_idx, best_objective = -1, -10000
 for idx in range(len(sentences)):
 if idx not in summary_idxes:
 redundancy = compute_redundancy(idx, summary_idxes)
 sentiment_difference = _sentiment_difference(idx, summary_idxes)
 # Maximizing centrality while minimizing redundancy and sentiment difference
 objective = centrality_scores[idx] - redundancy - sentiment_difference
 if objective > best_objective:
 best_idx = idx
 best_objective = objective

 if best_idx != -1:

```

```
objective = centr
if objective > be
```

```

best_objective = objective
if best_idx != -1:
 summary_idxes.append(best_idx)

Store the outputs
out_summary = ""
out_sentiment = ""
out_centralities = ""
summary_sentiments = []
summary_centralities = []
summary_redundancies = []

for idx in summary_idxes:
 summary_sentiments.append(sentiments[idx])
 summary_centralities.append(centrality_scores[idx])

compute redundancy of the specific sentence to all other sentences in the summary
summary_idxes.remove(idx)
redundancy = compute_redundancy(idx, summary_idxes)
summary_idxes.append(idx)
summary_redundancies.append(redundancy)

out_summary += sentences[idx] + "[SEP]"
out_sentiment += str(sentiments[idx]) + "[SEP]"
out_centralities += str(centrality_scores[idx]) + "[SEP]"

mean_sentiment = np.mean(summary_sentiments)
mean_centrality = np.mean(summary_centralities)
mean_redundancy = np.mean(summary_redundancies)

```

```
out_centralities += str(c)
```

```
mean_sentiment = np.mean(summary_sentiments)
mean_centrality = np.mean(summary_centralities)
mean_redundancy = np.mean(summary_redundancies)

print(f"({mean_sentiment}, {mean_centrality}, {mean_redundancy}, {doc_sentiment}")

return out_summary, out_sentiment, out_centralities, mean_sentiment, mean_centrality, mean_redundancy, doc
```

### 3.5.3 Result & Evaluation

Due to lack of gold summaries, we use **three metrics (Centrality, Sentiment Differences, and Redundancy)** of extracting the top 5 central sentences

- **Centrality:** Average centrality of summary sentences  $\rightarrow$  (Higher the better)
- **Sentiment Difference:** |mean of original text - mean of summary [sentiment]|  $\rightarrow$  (Lower the better)
- **Redundancy:** Average cosine similarity among summary sentences  $\rightarrow$  (Lower the better)

```
In [204]: # Read result CSVs
summary_baseline = pd.read_csv('..output/summ-baseline-out.csv').iloc[:, 1:] # baseline model
summary_balanced = pd.read_csv('..output/summ-balance-out.csv').iloc[:, 1:] # balanced model
```

We define an evaluation function to compute the above 3 metrics

```
In [210]: def evaluation(name, stats_df):
 stats_df.replace(np.nan, np.nan, inplace=True) # replace 'inf' by 'NaN'
 df = df.dropna(axis=1, how='any', inplace=False)
```

We define an evaluation function

```
[210.] def evaluation(name, stats_df):
 stats_df.replace(np.inf, np.nan, inplace=True) # replace 'inf' by 'NaN'
 stats_df.dropna(inplace=True, how='any', inplace=True) # then drop all NaNs
 mean_df = stats_df.describe().loc['mean', :]

 print(name + " model metrics:")
 print("Centrality: ", mean_df['Mean centrality'])
 print("Sentiment Difference: ", np.abs(mean_df['Mean_sentiment'] - mean_df['doc_sentiment']))
 print("Redundancy: ", mean_df['Mean_redundancy'])

 pass

In [231.] evaluation("Baseline", summ_baseline.iloc[:, -4:])

Baseline model metrics:
Centrality: 1.457832593624464
Sentiment Difference: 0.034905608279972744
Redundancy: 0.659224955303985

In [233.] evaluation("Balanced", summ_balanced.iloc[:, -4:])

Balanced model metrics:
Centrality: 1.407710446431976
Sentiment Difference: 0.006182471345803298
Redundancy: 0.5434062235968687

We summarize the metrics for two methods as follows:
```

Redundancy: 0.5434062235968687

|                      | Baseline | Balanced Summarization Algorithm |
|----------------------|----------|----------------------------------|
| Centrality           | 1.4578   | 1.4077                           |
| Sentiment Difference | 0.0349   | 0.0062                           |
| Redundancy           | 0.6592   | 0.5434                           |

We can observe that, though sacrificing minor centrality, our Improvement method (Balanced Summarization Algorithm) achieved **lower sentiment difference** and **lower redundancy** than baseline method, generally a good result.

---

## Part 4) Summary & Future plans

In this module, we first conduct some modifications on original datasets based on findings from exploratory data analysis and prepare it into a less sparse form for model part. Then we explore multiple recommendation models with optimizations that produce impressive results. Finally, we employ NLP models to provide more detailed information alongside books recommended to users, playing a role of user interaction tool and effective information extraction.

In the next module, we shall consider the following advancement and plans:

- Ensemble model:** Currently we have considered several recommendation models and a hybrid one to combine them. Next we prepare to introduce some linear models that take book features & user features into consideration and then build a high-functional

In the next module, we shall consider the

- **Ensemble model:** Currently we have considered several recommendation models and a hybrid one to combine them. Next we prepare to introduce some linear models that take book features & user features into consideration and then build a high-functional ensemble model as our final model;
- **NLP:** We will use NLP models to provide information overview of recommended book with our tags, summarization and sentiment scores. Then integrate models with interactive user interface to help users efficiently acquire information about books recommended to them
- **Interactive User Interface:** In the next module, we will present an interactive user interface to play with our recommendation system. It integrates recommendation models that produce to-read lists for users and use NLP models to provide information overview.