

INDENG 243 Project - Module 2

GoodReads: Book Analytics & Recommendation System

- @Group: 18
- @Member: Yanbo Wang, Yinyu Hu, Jiaming Xiong, Shichen Wu, Xinlin Huang, Wenxuan Yang
- @GitHub: <https://github.com/YanboWang2004/GoodReads-Recommendation-System>

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-; 'Fiction', '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 'Fiction', '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...	
	2671	15881	Harry Potter and the Chamber of Secrets	J.K. Rowling	fantasy 'Fiction', '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...	
	2869	15881	Harry Potter and the Chamber of Secrets	J.K. Rowling	fiction 'Fiction', '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...	

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 Two_Star 7776 non-null int64
 16 One_Star 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-cardinality-data-encoding/>

```
In [21]: Award_Num = df_info['Award'].apply(lambda x: 0 if x == [None] 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 its statistical indicators for better measurement:

- **Median**: measures the central tendency of data distribution
- **QR (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 13 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: object(1), int64(4), object(8)
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('int64').values
    data_df = pd.DataFrame(data=[uid] + list(data), index = ['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()

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

```
Out[39]:
```

	UId	AddedShelf	ToRead	added_median	ratings_median	reviews_median	to-read_median	added_IQR	ratings_IQR	reviews_IQR	to-read_IQR
0	10006486	13107	14731	7.0	0.0	0.0	4.0	1.000	2.000	0.0	0.0
1	10006486	2627	1670	7.0	0.0	0.0	4.0	1.000	0.000	0.0	0.0
2	10005151	124063	926493	24.0	7.0	1.0	7.0	8.000	7.000	4.0	0.0
3	10008056	36918	6637	12.0	4.0	0.0	2.0	6.000	6.000	1.000	0.0
4	10003977	27195	11902	22.0	6.0	0.0	9.0	17.75	5.000	1.0	1.0
...
7771	999713	22569	8889	3.0	1.0	0.0	2.0	3.00	1.000	0.0	0.0
7772	9984	244540	121383	101.0	28.0	1.0	59.0	39.000	12.000	2.0	2.0
7773	9994	192255	98182	58.0	10.0	1.0	35.0	18.000	7.000	1.0	1.0
7774	99955	35739	18080	5.0	1.0	0.0	3.0	2.50	1.000	0.0	0.0
7775	9999107	85198	32656	19.0	5.0	0.0	7.0	7.75	5.000	1.0	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: object(4), int64(6)
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 to see if we can in 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()
```

```
Out[119]:
```

	Reviewer	N_Review
	Amanda	31
	Luna	265
(shan) Littlebookcove	112	1
	The Polybrary	313
	A.	1163
	점스마	165
	여리고	167
	D r e y	71
	F r i e n d s	845
	e l l i e	182

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. Someday I will 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)
```

```
Out[128]: (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 Archaeology 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. Someday I will 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 [138]: from string import punctuation

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

    if not_punct == '':
        return document
    else:
        return not_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 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 = []
    review_group = []
    non_review_group = []

    flag = False
    while (temp_group != review_group) & (temp_non_group != non_review_group):
        print("Current user name: ", name, "id: ", name, "id: ")
        temp_group = df_reviews_copy[df_reviews_copy['Reviewer'] == name].groupby('N_Review').count()[0]
        temp_non_group = df_reviews_copy[df_reviews_copy['Reviewer'] == name].groupby('N_Review').count()[0]

        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()[0]['UId'])
            flag = True
            continue

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

            # if exists aggregated entries
            remove_group = {}
            remove_non_group = {}

            if review_group:
                for group v.s. independent sample
                for review_group in review_group:
                    if np.abs(review_num_1 - review_num_2) <= 100:
                        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()[0]['UId'])
                            df_reviews_copy.loc[index, 'N_Review'] = review_num_2 # change values
                            remove_group.append(review_num_1)
                        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_2)
                    print("-----")
                    review_group = [i for i in review_group if i not in set(remove_group)]
                    print(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 len(same_name_entries) > 1:
                                        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_non_group.append(non_group_1)
                                    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_non_group.append(non_group_2)
                                print("-----")
                                non_review_group = [i for i in non_review_group if i not in set(remove_non_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

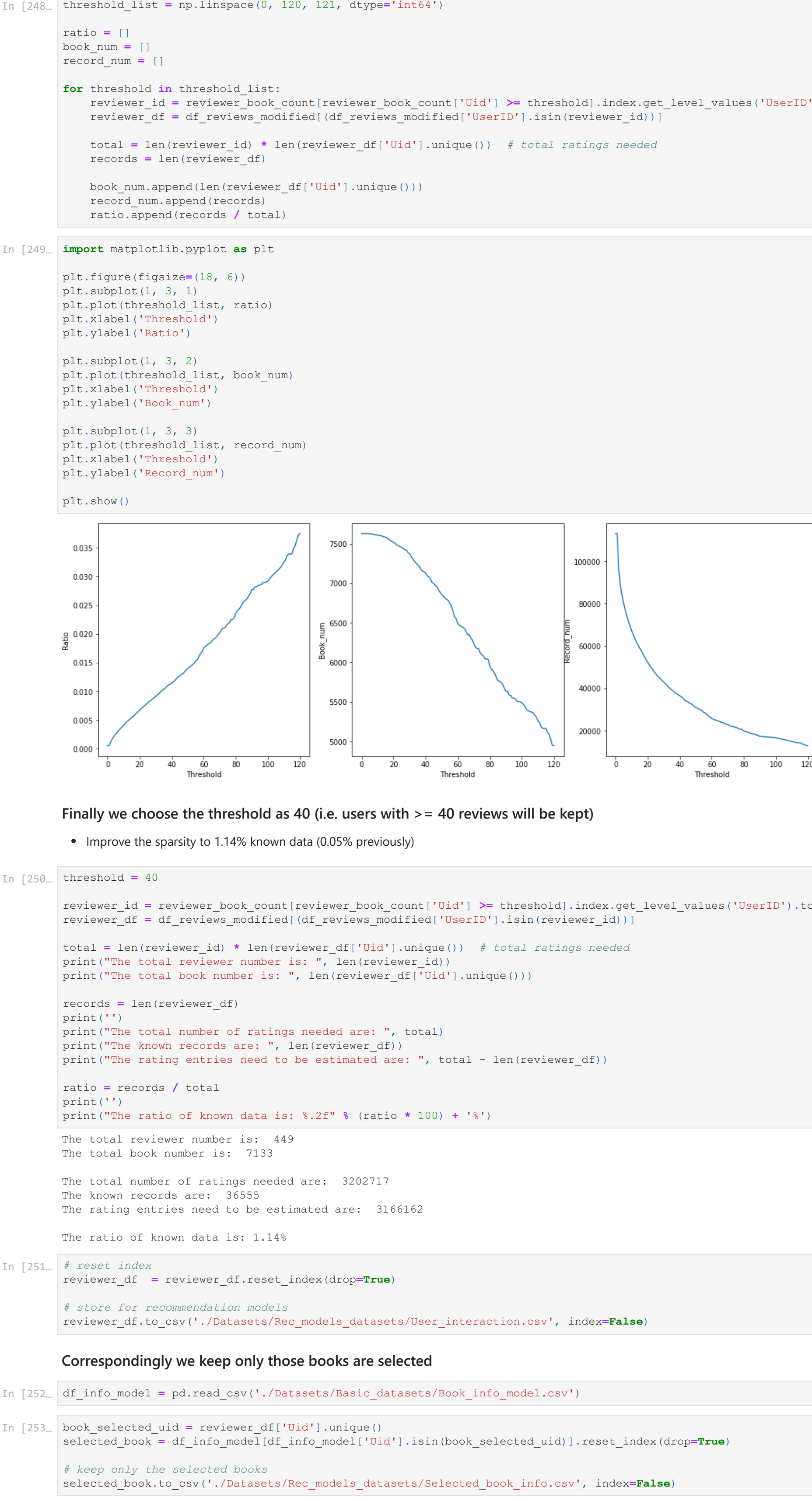
- Minor mistake in regular expressions causing the `N_Follower` data is wrong: only matched the first digit
- We use the previous combination result, taking the highest `N_Follower` number of the same user to replace all entries

```
In [29]: for name in names:
    same_name_entries = df_reviews_copy[df_reviews_copy['Reviewer'] == name]
    if len(same_name_entries) > 1:
        temp = same_name_entries.groupby(by=['N_Review']).count()[0]['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

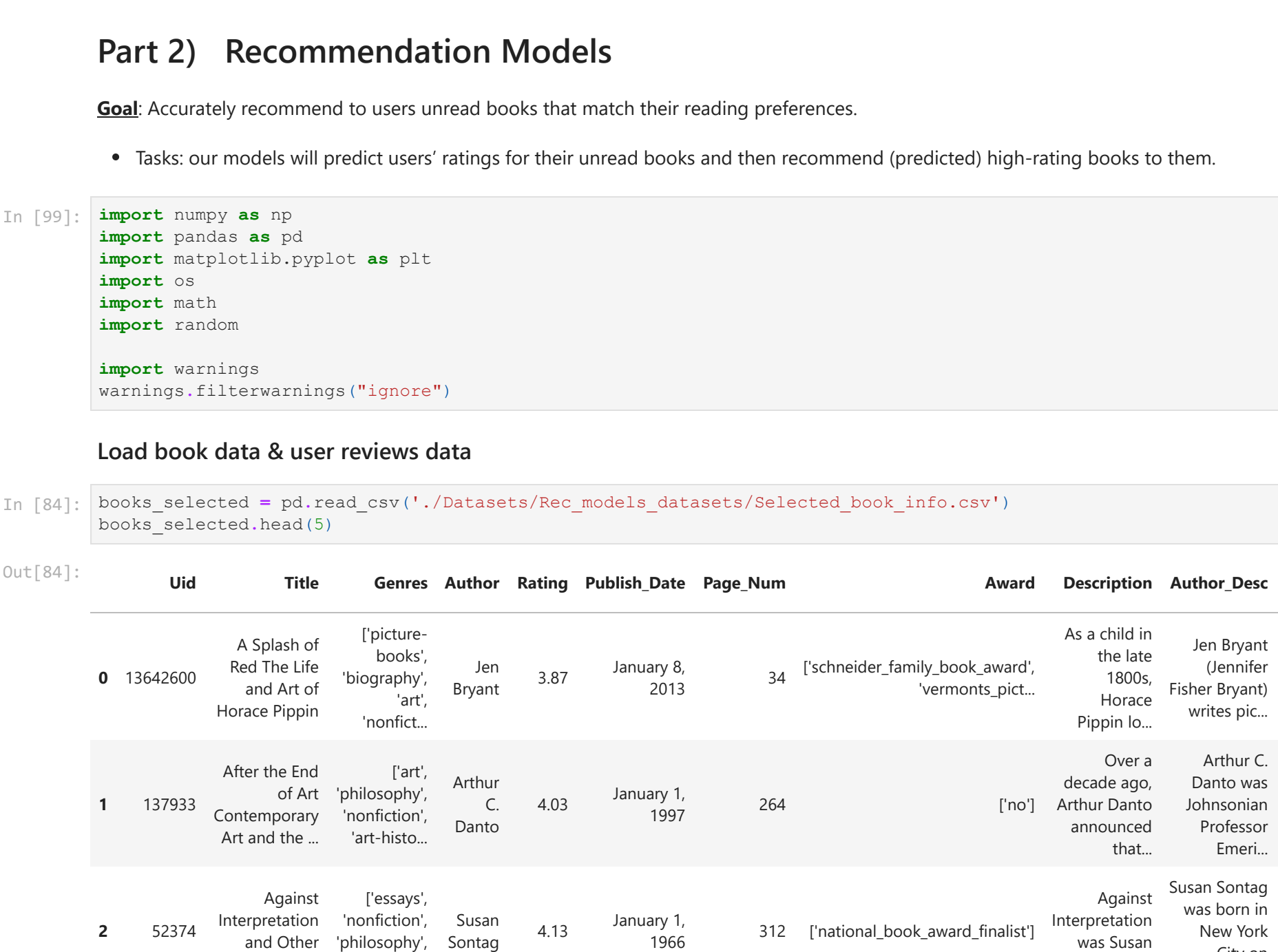
```
In [222]: df_reviews_copy.loc[:, 'UserID'] = df_reviews_copy.groupby(['Reviewer', 'N_Review']).nngroup()
```

```
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.columns = cols</
```

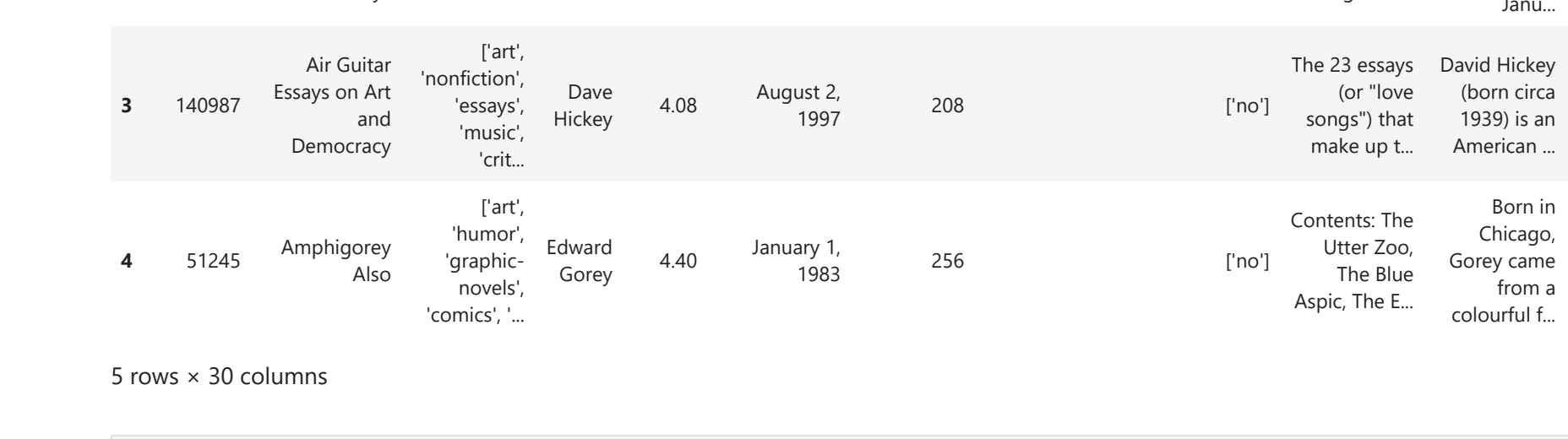



Finally we choose the threshold as 40 (i.e. users with >= 40 reviews will be kept)

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



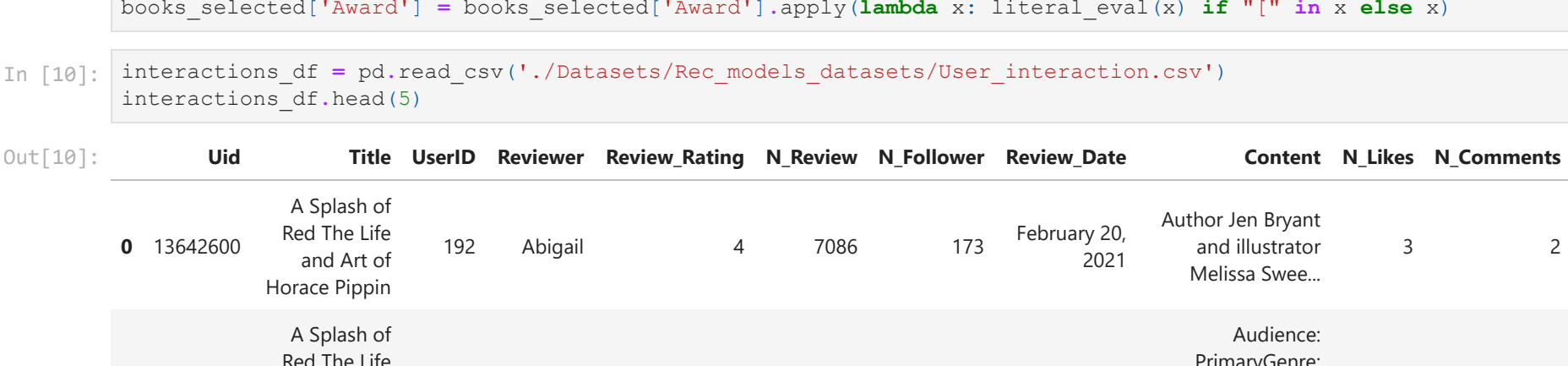
Correspondingly we keep only those books are selected



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 Fisher Bryant) writes pic...

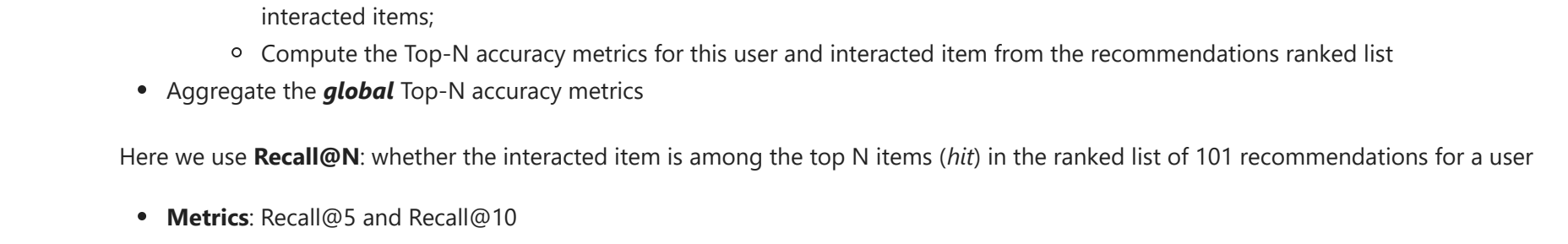
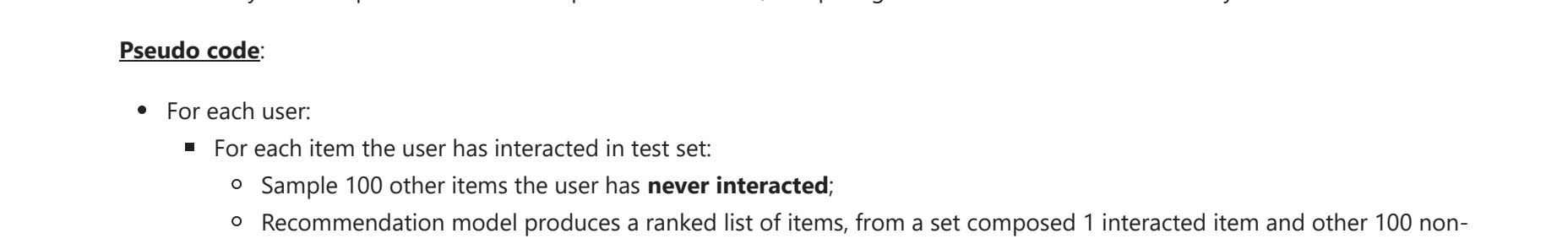
1	1379933	After the End of Art Contemporary Art and the ...	['art', 'philosophy', 'nonfiction', 'art-histo...']	Arthur C. Danto	4.03	January 1, 1997	264	['no']	Over a decade ago, Arthur C. Danto announced that...	Arthur C. Danto was a Professor Emerit...
---	---------	---	---	-----------------	------	-----------------	-----	--------	--	---

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 Essays on Art and Democracy	['art', 'nonfiction', 'essays', 'music', 'crit...']	Dave Hickey	4.08	August 2, 1997	208	['no']	The 23 essays (for "How songs") that make up L...	Dave Hickey (for "How 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



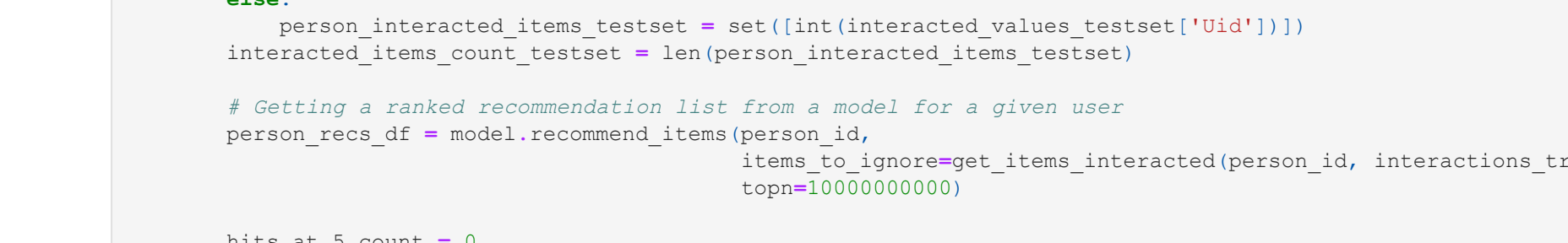
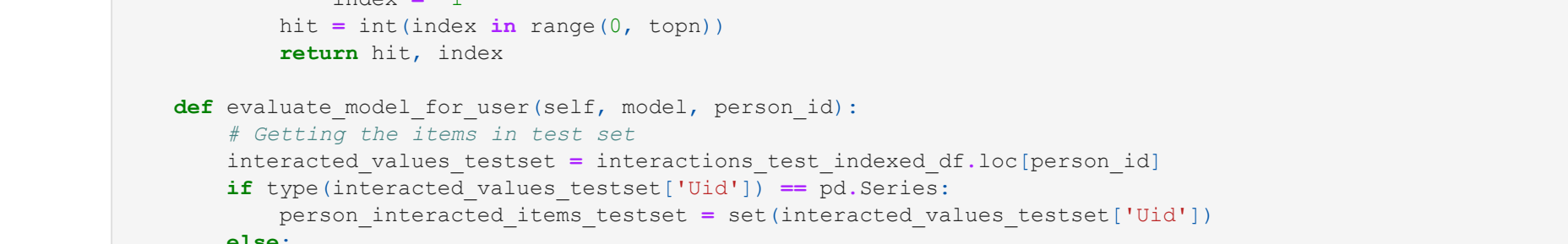
1	13642600	A Splash of Red The Life and Art of Horace Pippin	17464	Michelle	4	301	38	July 10, 2013	Audience: Primary/Genre: Non-Fiction/Informatio...	3	0
---	----------	---	-------	----------	---	-----	----	---------------	--	---	---

2	1379933	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	1379933	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:

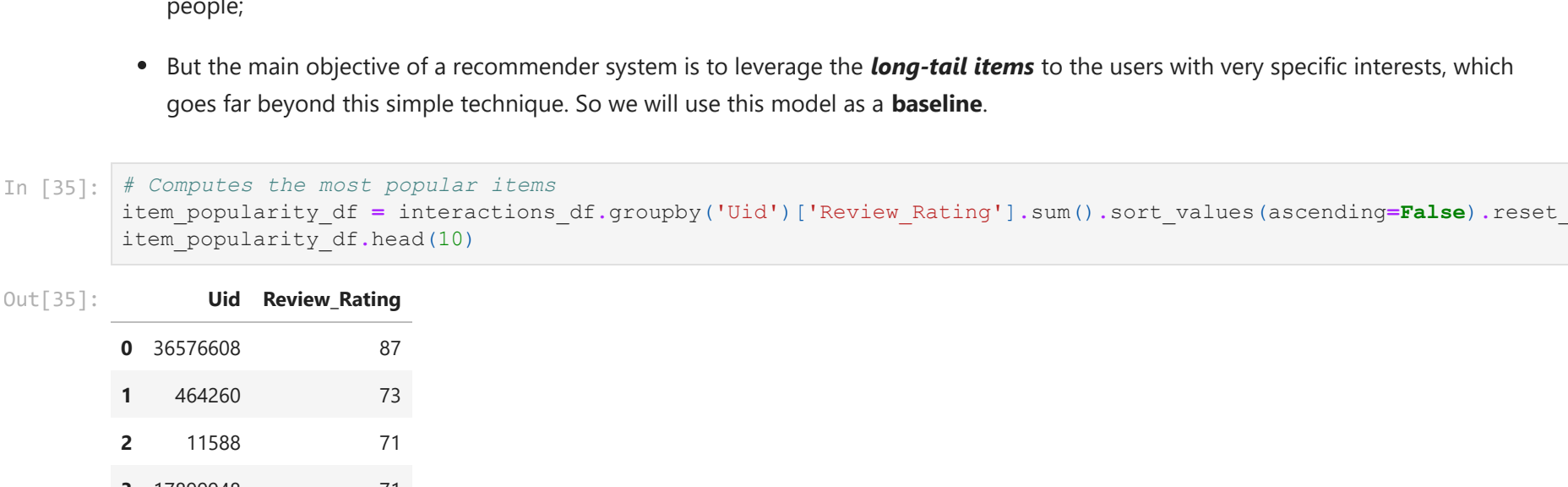
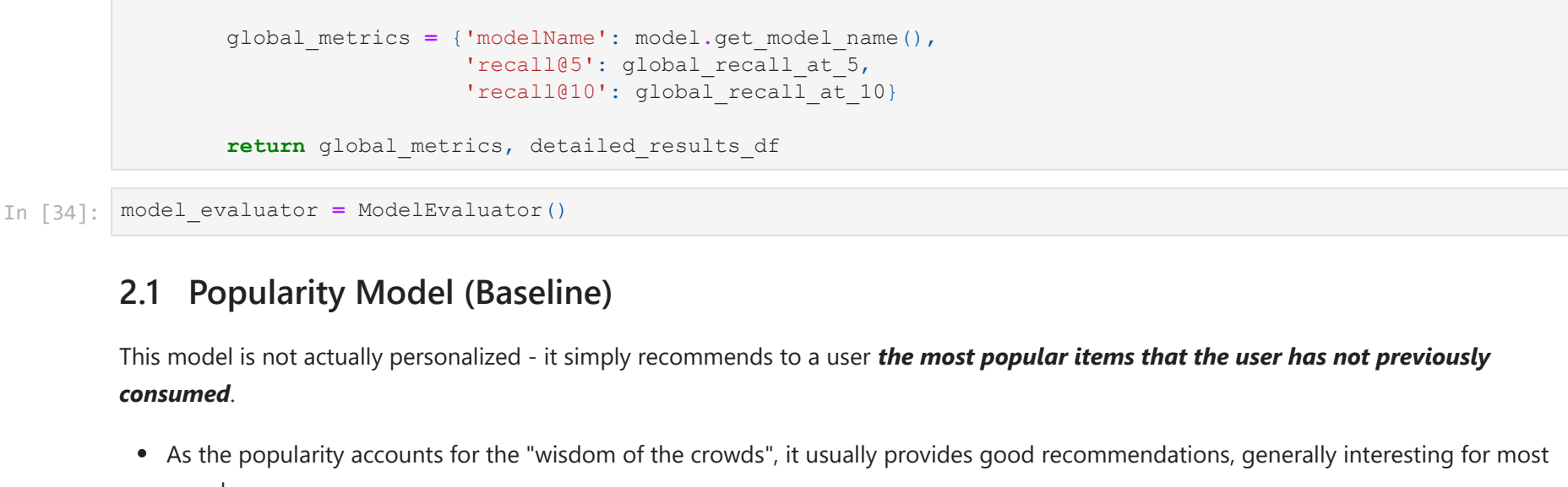
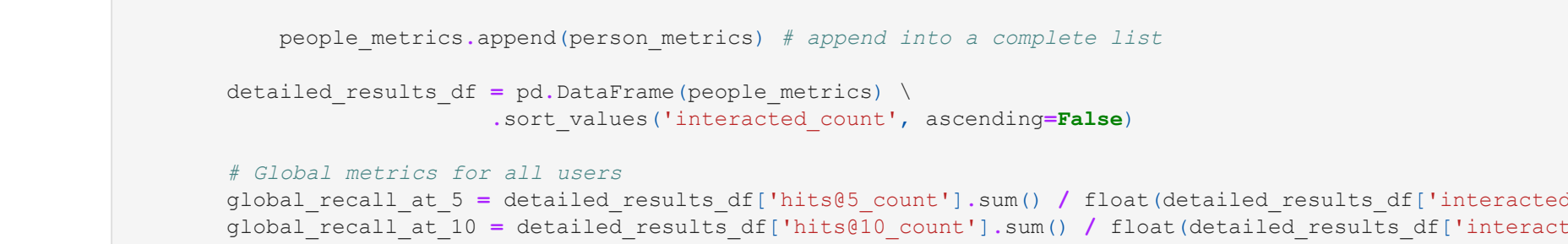
- Accuracy of the top recommendations provided to a user, comparing to the items the user has actually interacted in test set

Pseudo code:

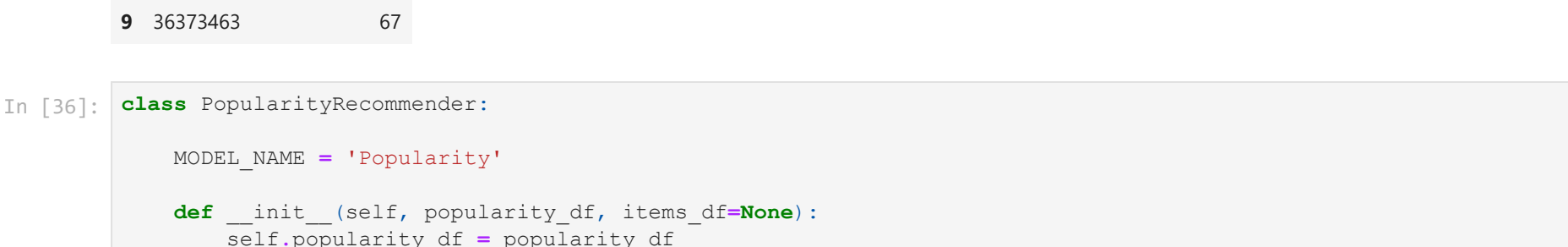
- For each user:
 - For each item the user has interacted in test set:
 - Sample 100 other items the user has **not** interacted with
 - Recommendation model produces a ranked list of items, from a set composed 1 interacted item and other 100 non-interacted items
 - Compute the Top-N accuracy metrics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



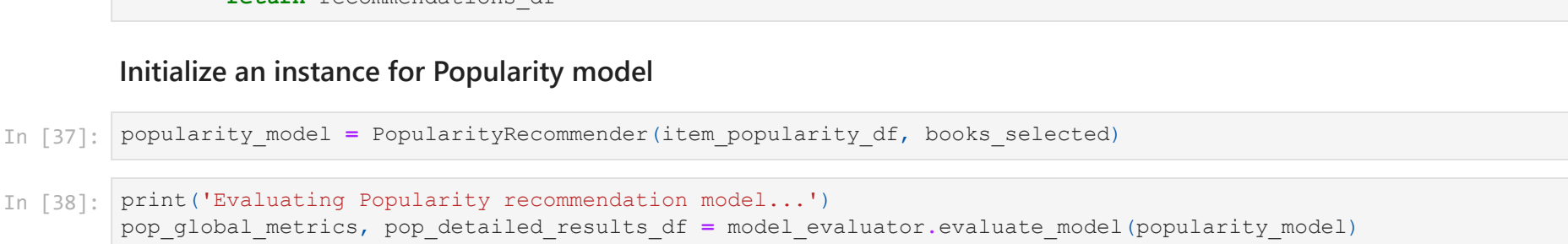
Global metrics for all users



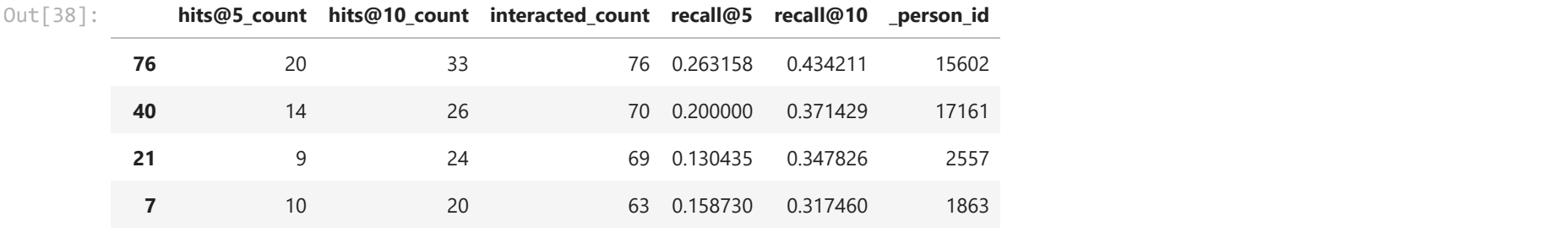
2.1 Popularity Model (Baseline)

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

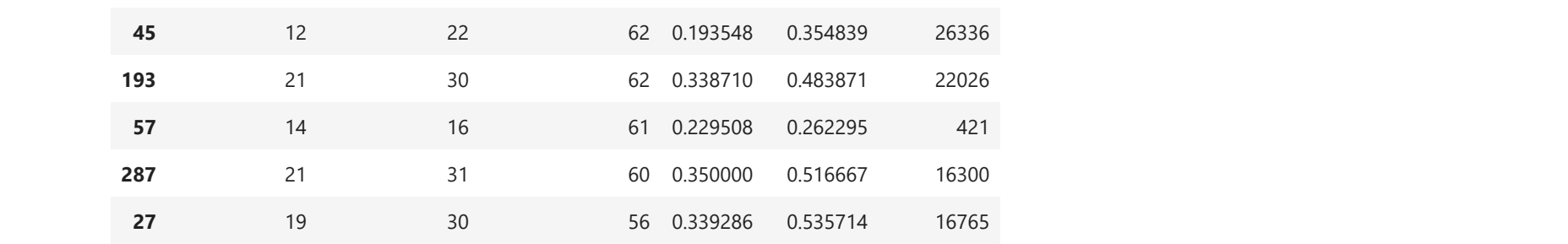
- As the popularity accounts for the 'wisdom of the crowds', it usually provides good recommendations, generally interesting for most people:
- But the main objective of a recommender system is to leverage the **long-tail items** to the users with very specific interests, which goes far beyond this simple 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	96358	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

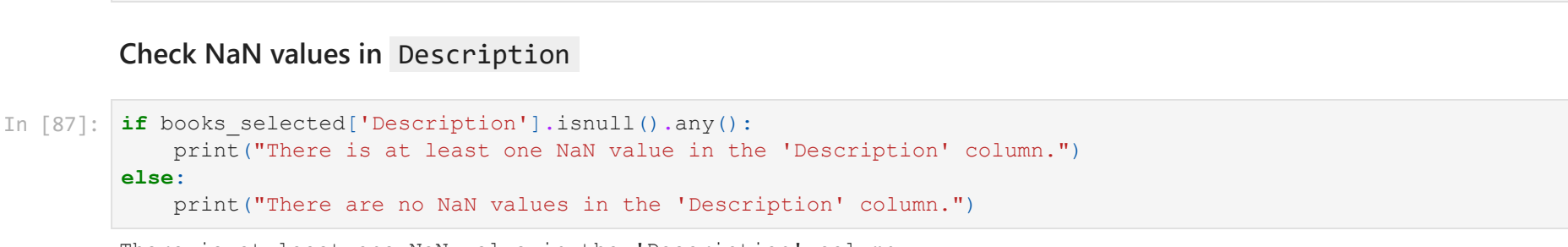
2.2 Content-based Filtering

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

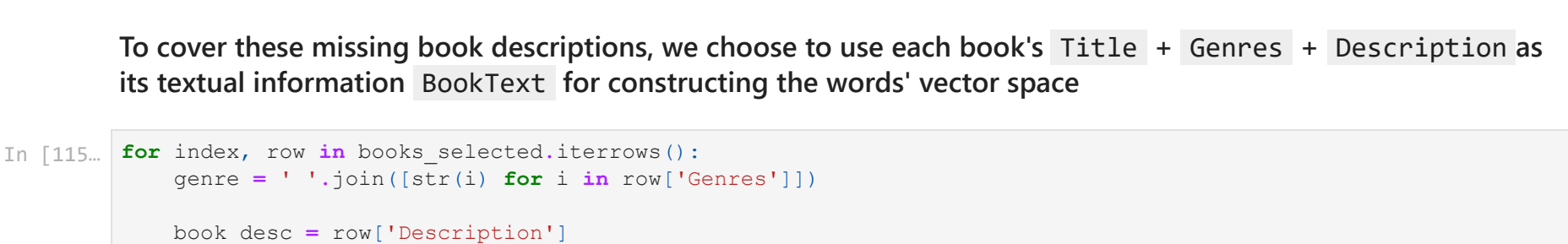
- It builds **users' profiles (tastes)** depending only on the user previous choices, making this method robust to avoid the cold-start problem

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

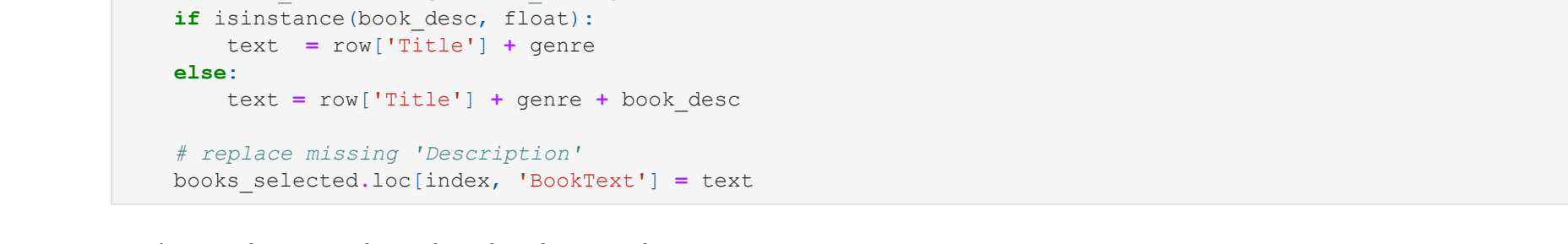
- This technique converts unstructured text into a vector structure, where each word is represented by a position in the vector, and the value measures **how relevant a given word is for a book**;
- As all items will be represented in the same **Vector Space Model**, we will compute **cosine similarity**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



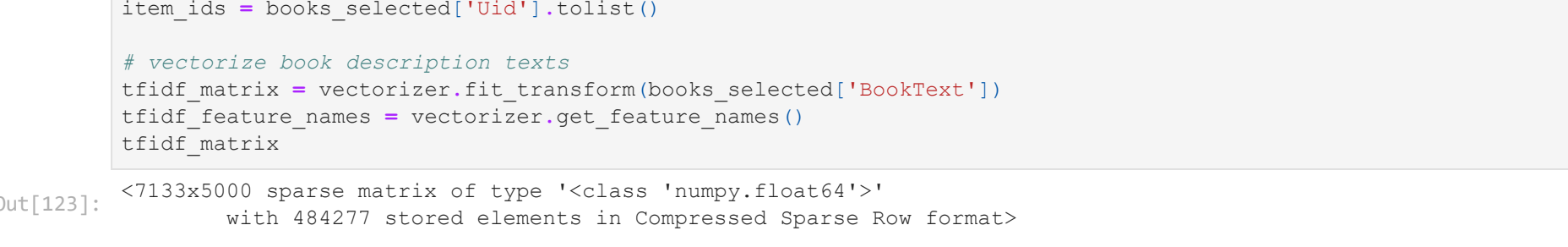
Check NaN values in Description



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

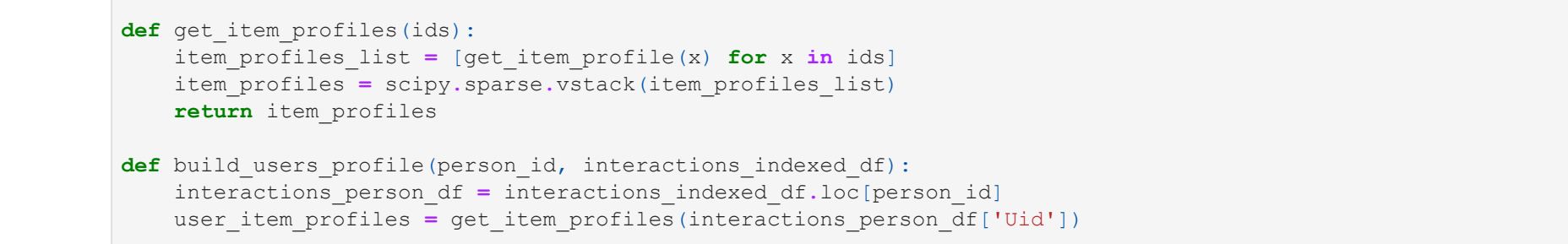
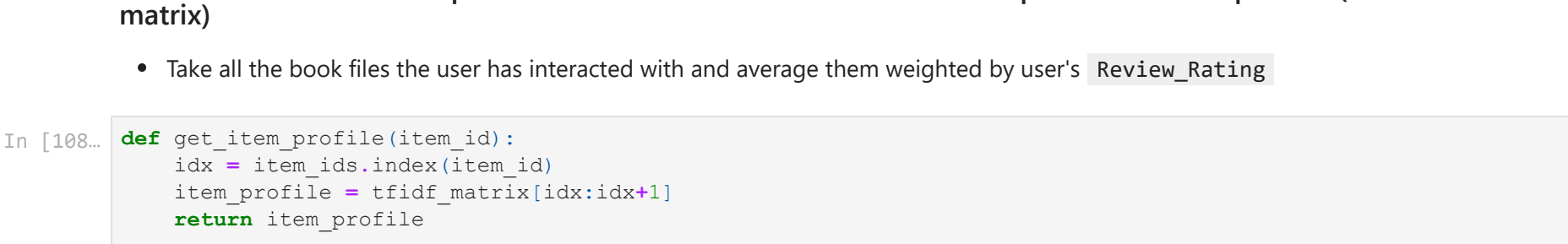


Train word vectors based on book texts data



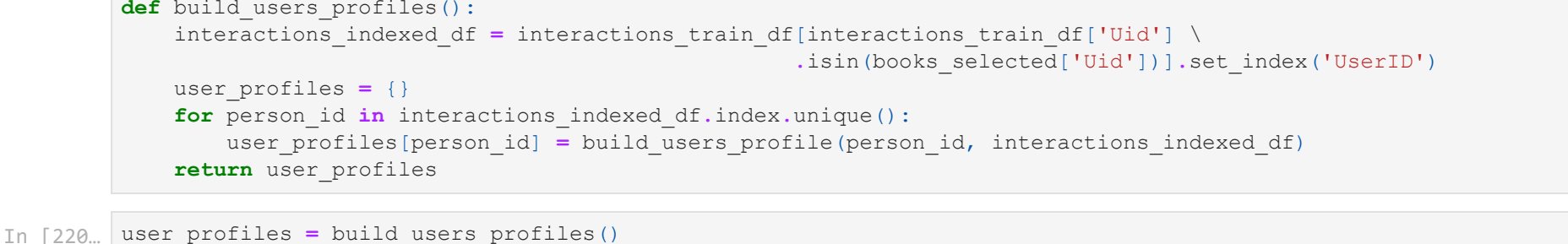
Obtain the word vector representation of a book in the same feature space as the user profiles (i.e. TF-IDF matrix)

- 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	0.112380
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)**.

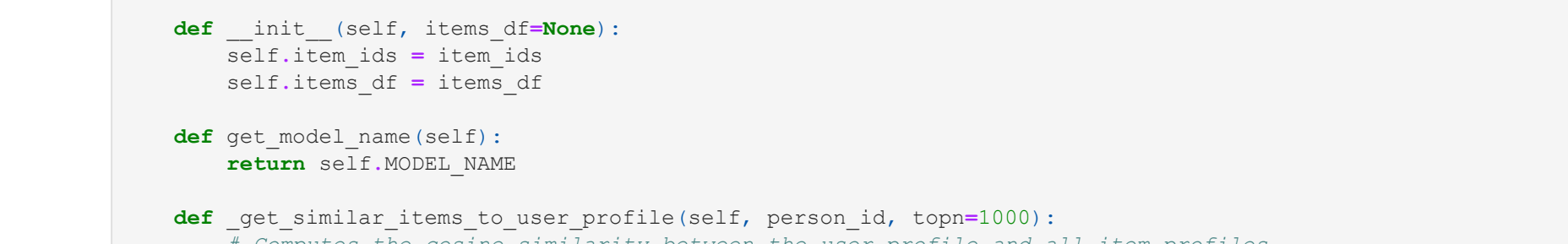
Singular Value Decomposition:

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

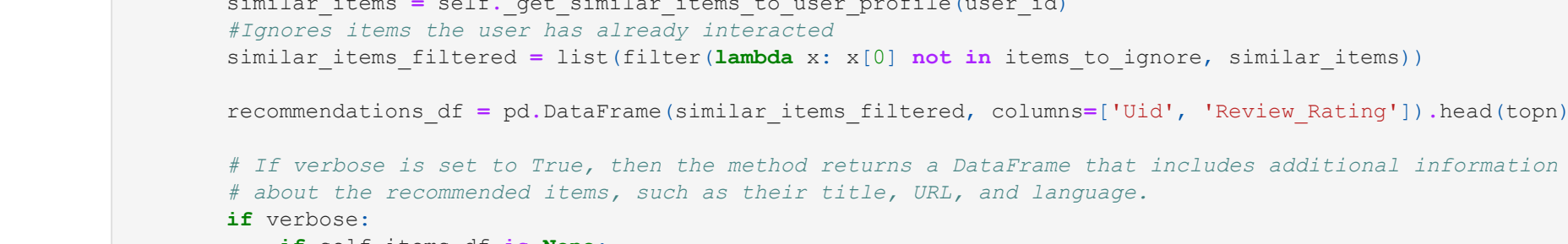
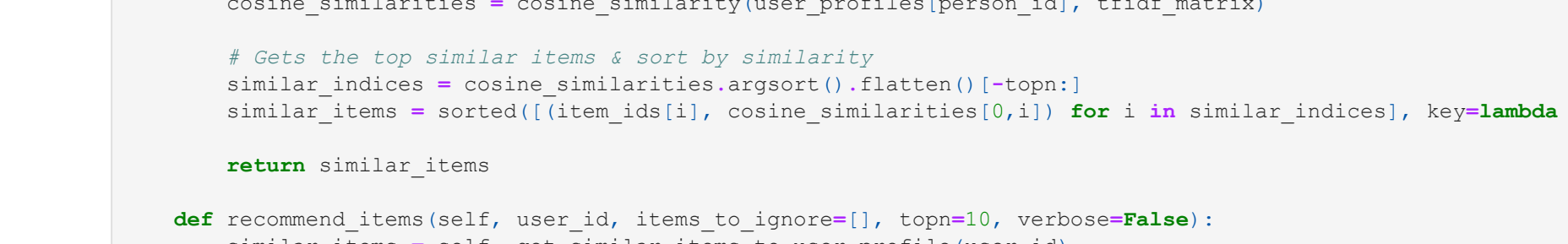
However, we can choose the **number of latent factors k** to factor the user-item matrix and approximate original matrix using the following formula:

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

- 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




```
154: class CFRecommender:
155:     MODEL_NAME = 'Collaborative Filtering'
156:
157:     def __init__(self, cf_predictions_df, items_df=None):
158:         self.cf_predictions_df = cf_predictions_df
159:         self.items_df = items_df
160:
161:     def get_model(self):
162:         return self.MODEL_NAME
163:
164:     def recommend_items(self, user_id, items_to_ignore=[], topn=10, verbose=False):
165:         # Get and sort the user's predictions
166:         sorted_user_predictions = self.cf_predictions_df[user_id].sort_values(ascending=False) \
167:             .reset_index().rename(columns={'user_id': 'Review_Rating'})
168:
169:         # Recommend the highest predicted rating movies that the user hasn't seen yet.
170:         recommendations_df = sorted_user_predictions[sorted_user_predictions['id'].isin(items_to_ignore)] \
171:             .sort_values('Review_Rating', ascending=False) \
172:             .head(topn)
173:
174:         if verbose:
175:             if self.items_df is None:
176:                 raise Exception("Items_df is required in verbose mode")
177:             recommendations_df = recommendations_df.merge(self.items_df, how='left',
178:                 left_on='id',
179:                 right_on='id')[['Review_Rating', 'id', 'title']]
180:
181:         return recommendations_df
182:
183: In [156]: cf_recommender_model = CFRecommender(cf_preds_df, books_selected)
184:
185: In [157]: print('Evaluating Collaborative Filtering (SVD Matrix Factorization) model...')
186: cf_global_metrics, cf_detailed_results_df = model_evaluator.evaluate_model(cf_recommender_model)
187:
188: print('Global metrics:\n%s' % cf_global_metrics)
189: cf_detailed_results_df.head(10)
190: Evaluating Collaborative Filtering (SVD Matrix Factorization) model...
191: 448 users processed
192: Global metrics:
193: {'model_name': 'Collaborative Filtering', 'recall@5': 0.30488305239393514, 'recall@10': 0.41991519627957874}
194: Out[157]:
```

```
plt.xlabel('Number of factors (k) in SVD')
plt.ylabel('Metrics')
plt.grid()
plt.legend(loc='best')
plt.show()
```

CF Performance v.s. Number of factors

Recall

Number of factors (k) in SVD

Global Recall@5
Global Recall@10
Average performance

Best average

Find the optimal number of factors based on the average performance of metrics

```
index = average_metric.argmax()
k_optimal = num_factor[index]

print("The best number of factor k = ", k_optimal)
print("The corresponding Recall@5 = ", global_recall_5[index])
print("The corresponding Recall@10 = ", global_recall_10[index])

The best number of factor k = 23
The corresponding Recall@5 = 0.3036520311858843
The corresponding Recall@10 = 0.4275748871563398
```

maximum of two metrics, we can see no large differences with our chosen optimal k

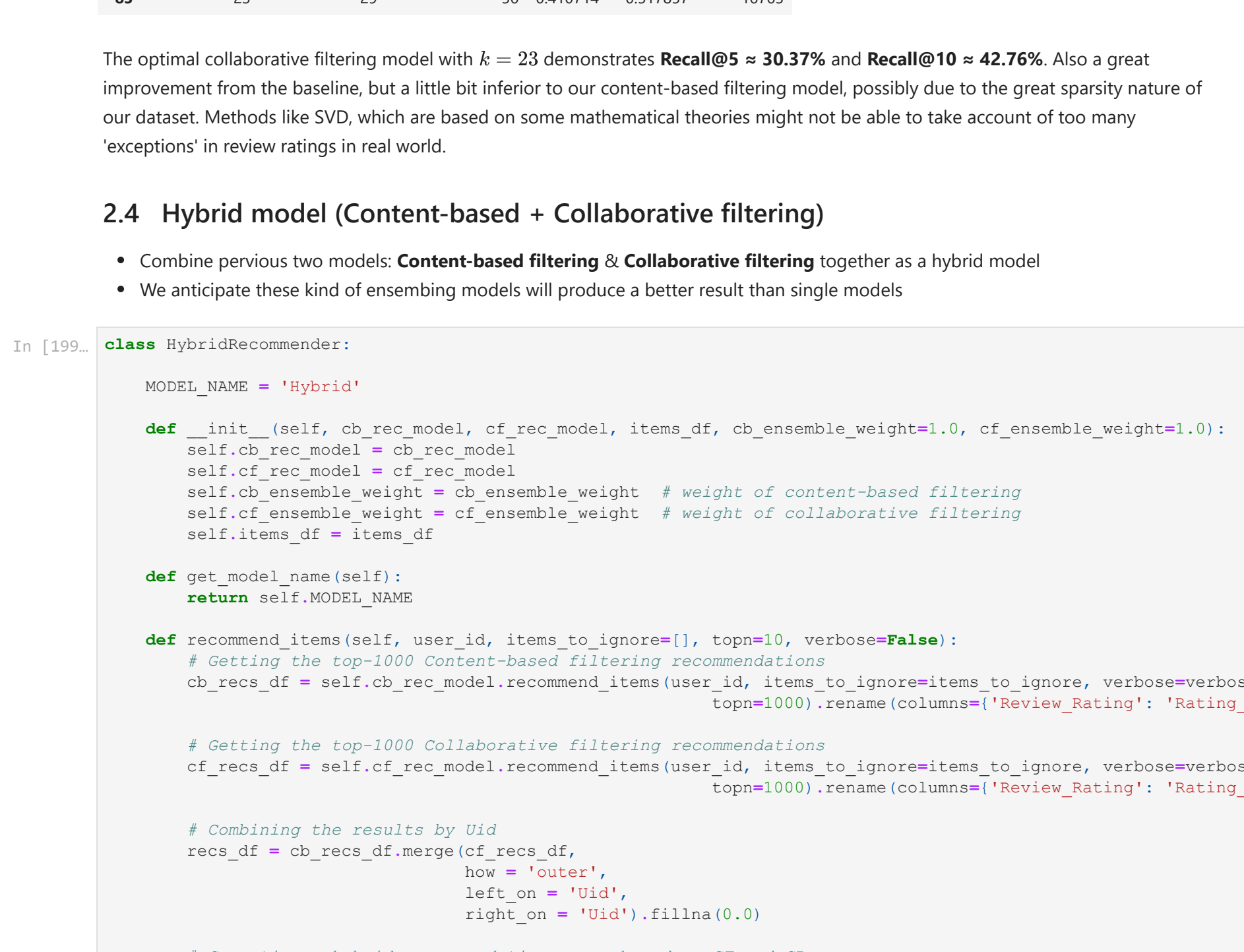
```
np.max(global_recall_5), np.max(global_recall_10)
```

(0.30830255778963206, 0.42976373026398576)

Evaluating the Collaborative Filtering model (SVD matrix factorization), we observe that we got **Recall@5 = 30.49%** and **Recall@10 = 41.99%**.

Next we iterate over a list of the number of factors to find the best number of factors *k*

```
In [ ]: num_factor = np.linspace(10, 100, 91, dtype='int64')
187:
188: global_recall_5 = []
189: global_recall_10 = []
190:
191: for n in num_factor:
192:     print("Number of factor is: ", n)
193:     U, sigma, Vt = svds(users_items_pivot_sparse_matrix, k = n)
194:     all_user_predicted_ratings = np.dot(np.dot(U, sigma), Vt)
195:     all_user_predicted_ratings_norm = (all_user_predicted_ratings - all_user_predicted_ratings.min()) / \
196:         (all_user_predicted_ratings.max() - all_user_predicted_ratings.min())
197:     cf_preds_df = pd.DataFrame(all_user_predicted_ratings_norm, columns = users_items_pivot_matrix_df.columns,
198:
199: # build model
200: cf_recommender_model = CFRecommender(cf_preds_df, books_selected)
201: cf_global_metrics, _ = model_evaluator.evaluate_model(cf_recommender_model)
202:
203: global_recall_5.append(cf_global_metrics['recall@5'])
204: global_recall_10.append(cf_global_metrics['recall@10'])
205:
206: Plot the metric performance v.s. number of factors (k)
```



Find the optimal number of factors based on the average performance of metrics

```
In [189]: index_optimal = np.argmax(k_optimal)
190: k_optimal = num_factor[index_optimal]
191:
192: print("The best number of factor k = ", k_optimal)
193: print("The corresponding Recall@5 = ", global_recall_5[index_optimal])
194: print("The corresponding Recall@10 = ", global_recall_10[index_optimal])
195:
196: The best number of factor k = 23
197: The corresponding Recall@5 = 0.3036520311858843
198: The corresponding Recall@10 = 0.4275748871563398
199:
200: In [191]: # Maximum of two metrics, we can see no large differences with our chosen optimal k
201: np.max(global_recall_5), np.max(global_recall_10)
202: Out[191]: (0.3036520311858843, 0.4275748871563398)
```

Re-train the model with optimal *k* = 23

```
In [198]: U, sigma, Vt = svds(users_items_pivot_sparse_matrix, k = k_optimal)
199: sigma = np.diag(sigma)
200:
201: all_user_predicted_ratings = np.dot(np.dot(U, sigma), Vt)
202: all_user_predicted_ratings_norm = (all_user_predicted_ratings - all_user_predicted_ratings.min()) / \
203:     (all_user_predicted_ratings.max() - all_user_predicted_ratings.min())
204: cf_preds_df = pd.DataFrame(all_user_predicted_ratings_norm, columns = users_items_pivot_matrix_df.columns, index
205:
206: # build model
207: cf_recommender_model = CFRecommender(cf_preds_df, books_selected)
208:
209: print('Evaluating Collaborative Filtering (SVD Matrix Factorization) model...')
210: cf_global_metrics, cf_detailed_results_df = model_evaluator.evaluate_model(cf_recommender_model)
211:
212: print('Global metrics:\n%s' % cf_global_metrics)
213: cf_detailed_results_df.head(10)
214: Evaluating Collaborative Filtering (SVD Matrix Factorization) model...
215: 448 users processed
216: Global metrics:
217: {'model_name': 'Collaborative Filtering', 'recall@5': 0.3036520311858843, 'recall@10': 0.4275748871563398}
218: Out[198]:
```

Metric	Recall@N	Recall@5	Recall@10
Recall@5	0.125	0.364	0.304
Recall@10	0.226	0.487	0.428

Though our models demonstrate some impressive results, but we should also pay attention to some potential limitations:

- Here we use **Recall@N** to measure the performance of models, but imagine for some 'extreme' user who is almost always being critical to books (i.e. always posts negative reviews), then our model only learns from the books **relevant** to him and it will be unlikely to recommend books that might not be the real favor of him.
- This can be extends to the *cold-start* problem for new users, if you only know one rating from this new user and that rating happens to be negative, the recommendation models might not work that well as desired;
- Therefore, currently we use users with at least 40 ratings, but we will try to test our trained models on users with less ratings (eliminated before) to simulate this situation and see how our models work and discuss potential ways to advance performance.

Part 3) Natural Language Processing (NLP) models

In this part, we consider some NLP models to prepare for the interaction tools that will be presented in the next module. The main functions are as follows:

- **Tag filtering system:** we want to extract **tags** for each book to enable users to filter books & comments by tags

The optimal collaborative filtering model with *k* = 23 demonstrates **Recall@5 = 30.37%** and **Recall@10 = 42.76%**. Also a great improvement from the baseline, but a little bit inferior to our content-based filtering model, possibly due to the great sparsity nature of our dataset. Methods like SVD, which are based on some mathematical theories might not be able to take account of too many 'exceptions' in review ratings in real world.

2.4 Hybrid model (Content-based + Collaborative filtering)

- Combine previous two models: **Content-based filtering** & **Collaborative filtering** together as a hybrid model
- We anticipate these kind of ensembling models will produce a better result than single models

```
In [199]: class HybridRecommender:
200:     MODEL_NAME = 'Hybrid'
201:
202:     def __init__(self, cb_rec_model, cf_rec_model, items_df, cb_ensemble_weight=1.0, cf_ensemble_weight=1.0):
203:         self.cb_rec_model = cb_rec_model
204:         self.cf_rec_model = cf_rec_model
205:         self.cb_ensemble_weight = cb_ensemble_weight # weight of content-based filtering
206:         self.cf_ensemble_weight = cf_ensemble_weight # weight of collaborative filtering
207:         self.items_df = items_df
208:
209:     def get_model_name(self):
210:         return self.MODEL_NAME
211:
212:     def recommend_items(self, user_id, items_to_ignore=[], topn=10, verbose=False):
213:         # Getting the top-1000 Content-based filtering recommendations
214:         cb_rec_df = self.cb_rec_model.recommend_items(user_id, items_to_ignore=items_to_ignore, verbose=verbose,
215:             topn=1000).rename(columns={'Review_Rating': 'Rating'})
216:         # Getting the top-1000 Collaborative filtering recommendations
217:         cf_rec_df = self.cf_rec_model.recommend_items(user_id, items_to_ignore=items_to_ignore, verbose=verbose,
218:             topn=1000).rename(columns={'Review_Rating': 'Rating'})
219:
220:         # Combining the results by id
221:         recs_df = cb_rec_df.merge(cf_rec_df,
222:             how='outer',
223:             left_on='id',
224:             right_on='id').fillna(0.0)
225:
226:         # Computing a hybrid recommendation score based on CF and CB scores
227:         recs_df['Rating_Hybrid'] = (recs_df['Rating_CF'] * self.cb_ensemble_weight) \
228:             + (recs_df['Rating_CF'] * self.cf_ensemble_weight)
229:
230:         # Sorting recommendations by Hybrid score
231:         recommendations_df = recs_df.sort_values('Rating_Hybrid', ascending=False).head(topn)
232:
233:         if verbose:
234:             if self.items_df is None:
235:                 raise Exception("Items_df is required in verbose mode")
236:             recommendations_df = recommendations_df.merge(self.items_df, how='left',
237:                 left_on='id',
238:                 right_on='id')[['Rating_Hybrid', 'id', 'title']]
239:
240:         return recommendations_df
241:
242: We first simply use the balanced weights 1:1 for two models
```

- You can choose optimal weights by trial and error, or use linear regression to study the optimal weights

```
In [200]: Hybrid_recommender_model = HybridRecommender(content_based_recommender_model, cf_recommender_model, books_selected,
201:     cb_ensemble_weight=1.0, cf_ensemble_weight=1.0)
202:
203: print('Evaluating Hybrid model...')
204: Hybrid_global_metrics, hybrid_detailed_results_df = model_evaluator.evaluate_model(Hybrid_recommender_model)
205:
206: print('Global metrics:\n%s' % Hybrid_global_metrics)
207: hybrid_detailed_results_df.head(10)
208: Evaluating Hybrid model...
209: 448 users processed
210: Global metrics:
211: {'model_name': 'Hybrid', 'recall@5': 0.51538771850636, 'recall@10': 0.6209820817845562}
212: Out[201]:
```

```
reviews['Content'] = reviews['Content'].apply(str.lower) # lower case
reviews['Content'] = reviews['Content'].apply(remove_punctuation)
reviews['Content'] = reviews['Content'].apply(remove_digit)

reviews.head(1)
```

	UId	Content
0	13642600	review later
1	13642600	we enjoyed this beautifully illustrated book...
2	13642600	featured in a grandma reads session exactly wh...
3	13642600	cute nonfiction picture book over the life of ...
4	13642600	a splash of red the life and art of horace p...l

Remove stopwords & frequent but less meaningful words

Here we define a **drop probability** (based from inverse frequency used in TD-IDF) to detect high-frequency words as follows:

$$\mathbb{P}_{Drop} = 1 - \frac{f}{frequency}$$

where t is a very small number and we will drop the words if $\mathbb{P}_{Drop} >= threshold$ to recude model noise and accelerate training

```
import nltk
from nltk.corpus import stopwords

# nltk.download('stopwords')
stop_words = set(stopwords.words('english'))
stop_words.update(['book', 'one', 'but', 'much', 'story', 'say', 'im', 'ive'])

from collections import Counter
from itertools import chain

def remove_stopwords(review):
    return [(word for word in simple_preprocess(str(review)) if word not in stop_words) for review in reviews]

# some words might appear too frequently but less meaningful, we also remove them
def remove_frequency_stop_words(reviews):
    t = 1e-5 # t-value
    threshold = 0.8 # threshold
    word_list = list(chain(*reviews)) # use chain() to combine all words

    # calculate word frequency
```

Comparing these models' performance



Though our models demonstrate some impressive results, but we should also pay attention to some potential limitations:

- Here we use **Recall@N** to measure the performance of models, but imagine for some 'extreme' user who is almost always being critical to books (i.e. always posts negative reviews), then our model only learns from the books **relevant** to him and it will be quite likely to recommend books that might not be the real favor of him.
- This can be extends to the **cold-start** problem for new users, if you only know one rating from this new user and that rating happens to be negative, the recommendation models might not work that well as desired.
- Therefore, currently we use users with at least 40 ratings, but we will try to test our trained models on users with less ratings (eliminated before) to simulate this situation and see how our models work and discuss potential ways to advance performance.

Part 3) Natural Language Processing (NLP) models

In this part, we consider some NLP models to prepare for the interaction tools that will be presented in the next module. The main functions are as follows:

- Tag filtering system**: we want to extract **tags** for each book to enable users to filter books & comments by tags
 - Topic Modeling** (Section 3.2): LDA method
 - Keywords Extraction** (Section 3.3): TextRank
 - Keyphrases Extraction** (Section 3.4): Yakei, keyBERT, Rake, etc.
- Extractive Summarization** (Section 3.5): Extract top review sentences to serve as a **summary** of all reviews of each book
 - Baseline approach
 - Balanced summarization

```
In [1]: import pandas as pd
2: import numpy as np
3: import re
4: import warnings
5: warnings.filterwarnings("ignore")
6:
7: 3.1 Texts Pre-processing
8:
9: Before we start the NLP models, we first conduct some pre-processing on our review content
```

	Uid	Title	Reviewer	N_Review	N_Follower	Review_Rating	Review_Date	Content	N_Likes	N_Comments
0	13642600	A Splash of Red The Life and Art of Horace Pippin	Chaplain Walle	380	18	4	January 22, 2023	review later	22	0
1	13642600	A Splash of Red The Life and Art of Horace Pippin	Lita	750	131	5	November 7, 2014	We enjoyed this beautifully illustrated book a...	11	2
2	13642600	A Splash of Red The Life and Art of Horace Pippin	Moonkist	1922	205	3	January 30, 2022	Featured in a grandma reads session. Exactly wh...	7	0
3	13642600	A Splash of Red The Life and Art of Horace Pippin	Cathy Newton	600	127	4	June 20, 2020	Cute nonfiction picture book over the life of...	6	0
4	13642600	A Splash of Red The Life and Art of Horace Pippin	Agne	744	57	3	August 16, 2017	A Splash of Red: The Life and Art of Horace Pi...	6	9

Here for NLP models, we are more interested in the textual reviews left by users: 'Content'

```
In [3]: reviews = df_reviews[['uid', 'Content']].copy()
4: reviews.head(5)
5: Out[3]:
```

uid	Content
0	13642600 review later
1	13642600 We enjoyed this beautifully illustrated book a...
2	13642600 Featured in a grandma reads session. Exactly wh...
3	13642600 Cute nonfiction picture book over the life of...
4	13642600 A splash of red the life and art of horace pipp...

Clean the reviews content

- Lowercase
- Remove punctuation
- Remove digits
- Remove stopwords & frequent but less meaningful words

```
In [4]: from string import punctuation
5:
6: def remove_punctuation(document):
7:     no_punct_char = []
8:     for character in document:
9:         # replace sentence separation punctuation by empty char
10:         if character in ['.', ',', '?', '!', ':', ';']:
11:             no_punct_char.append(' ')
12:         continue
13:
14:     if character not in punctuation:
15:         no_punct_char.append(character)
16:
17:     no_punct = ''.join(no_punct_char)
18:
19:     return no_punct
20:
21: def remove_digit(document):
22:     no_digit = ''.join([character for character in document if not character.isdigit()])
23:
24:     return no_digit
25:
26: In [5]: reviews['Content'] = reviews['Content'].apply(str.lower) # lower case
27: reviews['Content'] = reviews['Content'].apply(remove_punctuation) # if you want to execute visualisation prep yourself
28: reviews['Content'] = reviews['Content'].apply(remove_digit)
29:
30: In [6]: reviews.head()
31: Out[6]:
```

uid	Content
0	13642600 review later
1	13642600 we enjoyed this beautifully illustrated book a...
2	13642600 featured in a grandma reads session exactly wh...
3	13642600 cute nonfiction picture book over the life of...
4	13642600 a splash of red the life and art of horace pipp...

Remove stopwords & frequent but less meaningful words

Here we define a **drop probability** idea from inverse frequency used in TD-IDF to detect high-frequency words as follows:

$$P_{Drop} = 1 - \frac{f}{\text{frequency}}$$

where *f* is a very small number and we will drop the word has $P_{Drop} \geq \text{threshold}$ to recude model noise and accelerate training

```
In [148]: import nltk
149: from nltk.corpus import stopwords
150:
151: stopwords = nltk.download('stopwords')
152: stop_words = set(stopwords.words('english'))
153: stop_words.update(['book', 'one', 'br', 'much', 'story', 'say', 'in', 'ive'])
154:
155: In [149]: from collections import Counter
156: from itertools import chain
157:
158: def remove_stopwords(reviews):
159:     return [word for word in simple_preprocess(str(review)) if word not in stop_words for review in reviews]
160:
161: # some words might appear too frequently but less meaningful, we also remove them
162: def remove_frequency_stop_words(reviews):
163:     t = 1e-5 # t-value
164:     threshold = 0.8 # threshold
165:     word_list = list(chain(*reviews)) # use chain() to combine all words
166:
167:     # calculate word frequency
168:     int_word_counts = Counter(word_list)
169:     total_count = len(word_list)
170:     word_freq = {w: c / total_count for w, c in int_word_counts.items()}
171:
172:     # calculate the drop probability for words
173:     prob_drop = {w: 1 - np.exp(-f / f) for w, f in word_freq.items()}
174:
175:     # if drop probability < threshold, we keep this word for training
176:     train_words = []
177:     for i in range(len(reviews)):
178:         for w in reviews[i]:
179:             if prob_drop[w] < threshold:
180:                 train_words.append(w)
181:
182:     return train_words
183:
184: We first need to tokenize the reviews to remove stopwordsgensim to implement LDA
239: • Since we have to set the number of topics to train the LDA model, we experimented different options from 2 to 10.gensim.corpora.Dictionary(): map each normalized words with an integer id as a dictionary
244: • Then we use doc2bow() to convert document (a list of words) into the bag-of-words format: list of (token_id, token_count) 2-tuples.verbose=False because the model might need a long time for training, if you want to train by yourself, just turn verbose=True
```