

Book Recommendation System

Objective

The objective of this project is to develop a book recommendation system that can accurately suggest books to users based on user rating histories (collaborative filtering), clustering, and book genre prediction (content-based). The system leverages Book Recommendation Dataset, which contains user demographics, book metadata, and ratings, aiming at recommending the most-likely interested books to the user and providing a personalized reading experience.

Motivation

With the rise of web services, recommender systems are becoming more and more important in our daily lives. Right now, book recommendations are often made by looking at simple factors like what is popular/trending. This recommendation approach does not dig deep into readers' preferences, often leading to generic suggestions that don't quite hit the mark. This project aims at designing a book recommendation system that digs deeper by using information about users' rating histories and book metadata to make more personalized suggestions. A successful system can help users discover books they really love but might not have found, targeting a better user satisfaction.

+ Code

+ Text

Explore the Data

The Book Recommendation Dataset

The Book Recommendation Dataset (<https://www.kaggle.com/datasets/arashnic/book-recommendation-dataset>) provides detailed data in user demographic, book information, and ratings. This provides a solid foundation for developing a personalized book recommendation system.

"The Book-Crossing dataset comprises 3 files:

Users

Contains the users. Note that user IDs (User-ID) have been anonymized and map to integers. Demographic data is provided (Location, Age) if available. Otherwise, these fields contain NULL-values.

Books

Books are identified by their respective ISBN. Invalid ISBNs have already been removed from the dataset. Moreover, some content-based information is given (Book-Title, Book-Author, Year-Of-Publication, Publisher), obtained from Amazon Web Services. Note that in case of several authors, only the first is provided. URLs linking to cover images are also given, appearing in three different flavors (Image-URL-S, Image-URL-M, Image-URL-L), i.e., small, medium, large. These URLs point to the Amazon web site.

Ratings

Contains the book rating information. Ratings (Book-Rating) are either explicit, expressed on a scale from 1-10 (higher values denoting higher appreciation), or implicit, expressed by 0."

```
1 !pip install swifter

Requirement already satisfied: swifter in /usr/local/lib/python3.10/dist-packages (1.4.0)
Requirement already satisfied: pandas>=1.0.0 in /usr/local/lib/python3.10/dist-packages (from swifter) (2.0.3)
Requirement already satisfied: psutil>=5.6.6 in /usr/local/lib/python3.10/dist-packages (from swifter) (5.9.5)
Requirement already satisfied: dask[dataframe]>=2.10.0 in /usr/local/lib/python3.10/dist-packages (from swifter) (2023.8.1)
Requirement already satisfied: tqdm>=4.33.0 in /usr/local/lib/python3.10/dist-packages (from swifter) (4.66.2)
Requirement already satisfied: click>=8.0 in /usr/local/lib/python3.10/dist-packages (from dask[dataframe]>=2.10.0->swifter) (8.1.7)
Requirement already satisfied: cloudpickle>=1.5.0 in /usr/local/lib/python3.10/dist-packages (from dask[dataframe]>=2.10.0->swifter) (2.2.1)
Requirement already satisfied: fsspec>=2021.09.0 in /usr/local/lib/python3.10/dist-packages (from dask[dataframe]>=2.10.0->swifter) (2023.6.0)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from dask[dataframe]>=2.10.0->swifter) (24.0)
Requirement already satisfied: partd>=1.2.0 in /usr/local/lib/python3.10/dist-packages (from dask[dataframe]>=2.10.0->swifter) (1.4.1)
Requirement already satisfied: pyyaml>=5.3.1 in /usr/local/lib/python3.10/dist-packages (from dask[dataframe]>=2.10.0->swifter) (6.0.1)
Requirement already satisfied: toolz>=0.10.0 in /usr/local/lib/python3.10/dist-packages (from dask[dataframe]>=2.10.0->swifter) (0.12.1)
Requirement already satisfied: importlib-metadata>=4.13.0 in /usr/local/lib/python3.10/dist-packages (from dask[dataframe]>=2.10.0->swifter) (7.1.0)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.0.0->swifter) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.0.0->swifter) (2023.4)
Requirement already satisfied: tzdata>=2022.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.0.0->swifter) (2024.1)
Requirement already satisfied: numpy>=1.21.0 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.0.0->swifter) (1.25.2)
Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.10/dist-packages (from importlib-metadata>=4.13.0->dask[dataframe]>=2.10.0->swifter) (3.18.1)
Requirement already satisfied: locket in /usr/local/lib/python3.10/dist-packages (from partd>=1.2.0->dask[dataframe]>=2.10.0->swifter) (1.0.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.2->pandas>=1.0.0->swifter) (1.16.0)
```

Import Libraries

```
1#@title Import Libraries
2
3import numpy as np
4import pandas as pd
5import matplotlib.pyplot as plt
6from sklearn.exceptions import ConvergenceWarning
7from sklearn.metrics.pairwise import cosine_similarity
8from sklearn.neighbors import NearestNeighbors
9from sklearn.cluster import KMeans
10import gensim
11import re
12from gensim.models import Word2Vec
13import swifter
14
15pd.set_option('display.float_format', lambda x: '%.2f' % x)
```

Mount Google Drive

```
1#@title Mount Google Drive
2
3from google.colab import drive
4drive.mount('/content/drive')

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
```

Read Data In

```
1#@title Read Data In
2
3books = pd.read_csv("/content/drive/MyDrive/book-recommendation-dataset/Books.csv")
4ratings = pd.read_csv("/content/drive/MyDrive/book-recommendation-dataset/Ratings.csv")
5users = pd.read_csv("/content/drive/MyDrive/book-recommendation-dataset/Users.csv")

<ipython-input-295-08e09dd2ffcc>:3: DtypeWarning: Columns (3) have mixed types. Specify dtype option on import or set low_memory=False.
books = pd.read_csv("/content/drive/MyDrive/book-recommendation-dataset/Books.csv")
```

Understand the Data

```
1# Check shape
2print(books.shape)
3print(users.shape)
4print(ratings.shape)

(271360, 8)
(278858, 3)
(1149780, 3)
```

```
1 books.head(5)
```

	ISBN	Book-Title	Book-Author	Year-Of-Publication	Publisher	Image-URL-S	
0	0195153448	Classical Mythology	Mark P. O. Morford	2002	Oxford University Press	http://images.amazon.com/images/P/0195153448.0...	http://images.amazon.co
1	0002005018	Clara Callan	Richard Bruce Wright	2001	HarperFlamingo Canada	http://images.amazon.com/images/P/0002005018.0...	http://images.amazon.co
2	0060973129	Decision in Normandy	Carlo D'Este	1991	HarperPerennial	http://images.amazon.com/images/P/0060973129.0...	http://images.amazon.co
3	0374157065	Flu: The Story of the Great Influenza Pandemic...	Gina Bari Kolata	1999	Farrar Straus Giroux	http://images.amazon.com/images/P/0374157065.0...	http://images.amazon.co
4	0393045218	The Mummies of Urumchi	E. J. W. Barber	1999	W. W. Norton & Company	http://images.amazon.com/images/P/0393045218.0...	http://images.amazon.co

```
1 ratings.head(5)
```

	User-ID	ISBN	Book-Rating	
0	276725	034545104X	0	
1	276726	0155061224	5	
2	276727	0446520802	0	
3	276729	052165615X	3	
4	276729	0521795028	6	

```
1 users.head(5)
```

	User-ID	Location	Age
0	1	nyc, new york, usa	NaN
1	2	stockton, california, usa	18.00
2	3	moscow, yukon territory, russia	NaN
3	4	porto, v.n.gaia, portugal	17.00
4	5	farnborough, hants, united kingdom	NaN

▼ Data Preprocessing

- 1. Drop the last several books columns containing image URLs which will not be useful for analysis.
- 2. Check and handle missing author and publisher value using mode.
- 3. Replace zero-value book ratings with non-zero ratings' median.
- 4. Replace NaN age value with ages' median.
- 5. Filter out invalid ages, non-positive years of publication.
- 6. Convert data type to integer (e.g. Age, Year-Of-Publication, Rating).
- 7. Merge dataframes.

▼ Preprocess Books

```
1 books.isnull().sum()
```

ISBN	0
Book-Title	0
Book-Author	2
Year-Of-Publication	0
Publisher	2
Image-URL-S	0
Image-URL-M	0
Image-URL-L	3
dtype:	int64

```
1 # Drop unwanted columns
2 books.drop(['Image-URL-S', 'Image-URL-M', 'Image-URL-L'], axis=1, inplace=True)
3
4 # Year of Publication: Change NaN value to 0
5 books['Year-Of-Publication'] = pd.to_numeric(books['Year-Of-Publication'], errors='coerce').fillna(0)
6 books['Year-Of-Publication'] = books['Year-Of-Publication'].astype(int)
7
8 # Notice that some years are apparently invalid (e.g. year = 2050)
9 books.describe()
```

	Year-Of-Publication
count	271360.00
mean	1959.74
std	258.08
min	0.00
25%	1989.00
50%	1995.00
75%	2000.00
max	2050.00

```
1 # Filter out invalid years of publication, set valid years of publications between 1900-2024)
2 books = books[books['Year-Of-Publication'] >= 1900]
3 books = books[books['Year-Of-Publication'] <= 2024]
4
5 # Handle null values
6 books['Book-Author'].fillna(books['Book-Author'].mode()[0], inplace=True)
7 books['Publisher'].fillna(books['Publisher'].mode()[0], inplace=True)
8 books.info()
```

<class 'pandas.core.frame.DataFrame'>				
Index: 266723 entries, 0 to 271359				
Data columns (total 5 columns):				
#	Column	Non-Null Count		Dtype
0	ISBN	266723 non-null		object
1	Book-Title	266723 non-null		object
2	Book-Author	266723 non-null		object
3	Year-Of-Publication	266723 non-null		int64
4	Publisher	266723 non-null		object
dtypes: int64(1), object(4)				
memory usage: 12.2+ MB				

```
1 books.head(5)
```

	ISBN	Book-Title	Book-Author	Year-Of-Publication	Publisher
0	0195153448	Classical Mythology	Mark P. O. Morford	2002	Oxford University Press
1	0002005018	Clara Callan	Richard Bruce Wright	2001	HarperFlamingo Canada
2	0060973129	Decision in Normandy	Carlo D'Este	1991	HarperPerennial
3	0374157065	Flu: The Story of the Great Influenza Pandemic...	Gina Bari Kolata	1999	Farrar Straus Giroux
4	0393045218	The Mummies of Urumchi	E. J. W. Barber	1999	W. W. Norton & Company

Preprocess Ratings

```
1 ratings.isnull().sum()
```

```
User-ID      0
ISBN         0
Book-Rating  0
dtype: int64
```

```
1 # Notice that more than half of the book ratings are 0 (which is invalid)
2 ratings.describe()
```

	User-ID	Book-Rating
count	1149780.00	1149780.00
mean	140386.40	2.87
std	80562.28	3.85
min	2.00	0.00
25%	70345.00	0.00
50%	141010.00	0.00
75%	211028.00	7.00
max	278854.00	10.00

```
1 # Replace 0 ratings with non-zero ratings' median
2 median_rating = ratings['Book-Rating'][ratings['Book-Rating'] != 0].median()
3 ratings['Book-Rating'].replace(0, median_rating, inplace=True)
4 ratings = ratings[ratings['Book-Rating'] > 0]
5 ratings.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1149780 entries, 0 to 1149779
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0   User-ID     1149780 non-null  int64
1   ISBN        1149780 non-null  object
2   Book-Rating 1149780 non-null  int64
dtypes: int64(2), object(1)
memory usage: 26.3+ MB
```

```
1 ratings.head(5)
```

	User-ID	ISBN	Book-Rating
0	276725	034545104X	8
1	276726	0155061224	5
2	276727	0446520802	8
3	276729	052165615X	3
4	276729	0521795028	6

Preprocess Users

```
1 users.isnull().sum()
```

```
User-ID      0
Location     0
Age      110762
dtype: int64
```

```
1 # Handle null age values, replace null values with median age
2 users['Age'].fillna(users['Age'].median(), inplace=True)
3 users['Age'] = users['Age'].astype(int)
4 users.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 278858 entries, 0 to 278857
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0   User-ID     278858 non-null  int64
1   Location    278858 non-null  object
```

```

2 Age      278858 non-null int64
dtypes: int64(2), object(1)
memory usage: 6.4+ MB

```

```

1 # Notice that some ages are apparently invalid (e.g. age = 0 and age = 244)
2 users.describe()

```

	User-ID	Age
count	278858.00	278858.00
mean	139429.50	33.66
std	80499.52	11.28
min	1.00	0.00
25%	69715.25	29.00
50%	139429.50	32.00
75%	209143.75	35.00
max	278858.00	244.00

```

1 # Filter out invalid ages (set valid ages between 5-100)
2 users = users[5 <= users['Age']]
3 users = users[users['Age'] <= 100]

```

```
1 users.head()
```

	User-ID	Location	Age
0	1	nyc, new york, usa	32
1	2	stockton, california, usa	18
2	3	moscow, yukon territory, russia	32
3	4	porto, v.n.gaia, portugal	17
4	5	farnborough, hants, united kingdom	32

Merge Dataframes

```

1 merged_df = pd.merge(books, ratings, on='ISBN')
2 df = pd.merge(merged_df, users, on="User-ID")
3 df.head()

```

	ISBN	Book-Title	Book-Author	Year-Of-Publication	Publisher	User-ID	Book-Rating	Location	Age
0	0195153448	Classical Mythology	Mark P. O. Morford	2002	Oxford University Press	2	8	stockton, california, usa	18
1	0002005018	Clara Callan	Richard Bruce Wright	2001	HarperFlamingo Canada	8	5	timmins, ontario, canada	32
2	0060973129	Decision in Normandy	Carlo D'Este	1991	HarperPerennial	8	8	timmins, ontario, canada	32

```
1 df.isnull().sum()
```

```

ISBN      0
Book-Title 0
Book-Author 0
Year-Of-Publication 0
Publisher  0
User-ID    0
Book-Rating 0
Location   0
Age        0
dtype: int64

```

```
1 df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1012553 entries, 0 to 1012552
Data columns (total 9 columns):
#   Column              Non-Null Count  Dtype
---  -
0   ISBN                1012553 non-null object
1   Book-Title          1012553 non-null object
2   Book-Author         1012553 non-null object
3   Year-Of-Publication 1012553 non-null int64
4   Publisher            1012553 non-null object
5   User-ID             1012553 non-null int64
6   Book-Rating         1012553 non-null int64
7   Location            1012553 non-null object
8   Age                 1012553 non-null int64
dtypes: int64(4), object(5)
memory usage: 69.5+ MB

```

```
1 df.describe()
```

	Year-Of-Publication	User-ID	Book-Rating	Age	
count	1012553.00	1012553.00	1012553.00	1012553.00	
mean	1995.30	140592.64	7.86	35.72	
std	7.30	80468.35	1.14	10.59	
min	1900.00	2.00	1.00	5.00	
25%	1992.00	70415.00	8.00	31.00	
50%	1997.00	141183.00	8.00	32.00	
75%	2001.00	211391.00	8.00	41.00	
max	2024.00	278854.00	10.00	100.00	

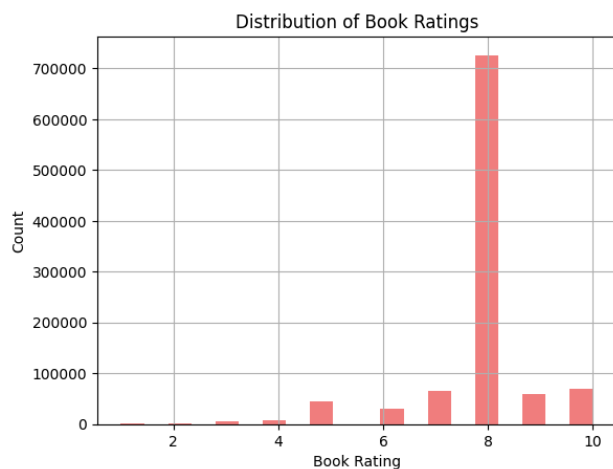
▼ Data Visualization

▼ Book Ratings

```

1 #@title Book Ratings
2 book_ratings = df['Book-Rating']
3 book_ratings.hist(bins=20, color='lightcoral')
4 plt.title('Distribution of Book Ratings')
5 plt.xlabel('Book Rating')
6 plt.ylabel('Count')
7 plt.show()

```

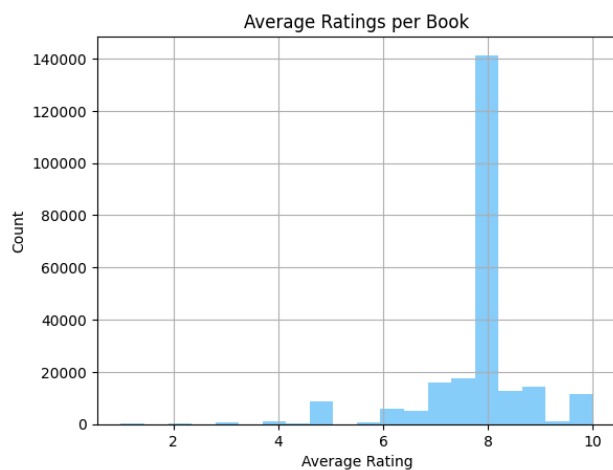


▼ Average Rating per Book

```

1 #@title Average Rating per Book
2 average_book_ratings = df.groupby('Book-Title')['Book-Rating'].mean()
3 average_book_ratings.hist(bins=20, color='lightskyblue')
4 plt.title('Average Ratings per Book')
5 plt.xlabel('Average Rating')
6 plt.ylabel('Count')
7 plt.show()

```

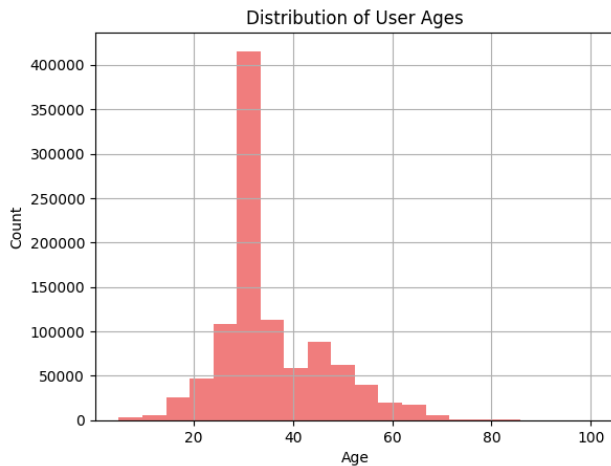


▼ User Ages

```

1#@title User Ages
2 ages = df['Age']
3 ages.hist(bins=20, color='lightcoral')
4 plt.title('Distribution of User Ages')
5 plt.xlabel('Age')
6 plt.ylabel('Count')
7 plt.show()

```

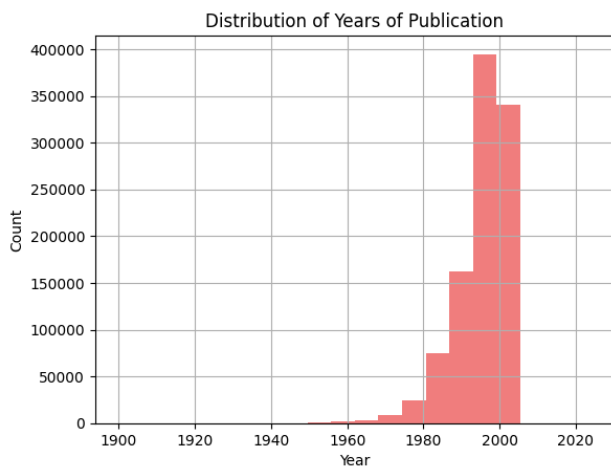


✓ Years of Publication

```

1#@title Years of Publication
2 publication_years = df['Year-Of-Publication']
3 publication_years.hist(bins=20, color='lightcoral')
4 plt.title('Distribution of Years of Publication')
5 plt.xlabel('Year')
6 plt.ylabel('Count')
7 plt.show()

```



✓ Filter Dataframe for Recommendation

After preprocessing, it can be observed that the DataFrame still remains considerably large, leading to a risk of exhausting RAM resources during analysis. Additionally, books with very few ratings and users who have provided only a limited number of ratings do not offer enough data points to effectively support the recommendation system. To address these issues, we should further refine our dataset to include only books that have received more than 100 ratings and users who have given more than 50 ratings to ensure we have more reliable data for analysis.

```

1# Filter books that have more than 100 ratings
2 book_counts = df['Book-Title'].value_counts()
3 books_over_100_ratings = book_counts[book_counts > 100].index
4
5# Filter users that have given more than 50 ratings
6 user_counts = df['User-ID'].value_counts()
7 users_over_50_ratings = user_counts[user_counts > 50].index
8
9# Filter the DataFrame that contains only selected books and users
10 filtered_df = df[(df['Book-Title'].isin(books_over_100_ratings)) & (df['User-ID'].isin(users_over_50_ratings))]

```

```
1 filtered_df.describe()
```

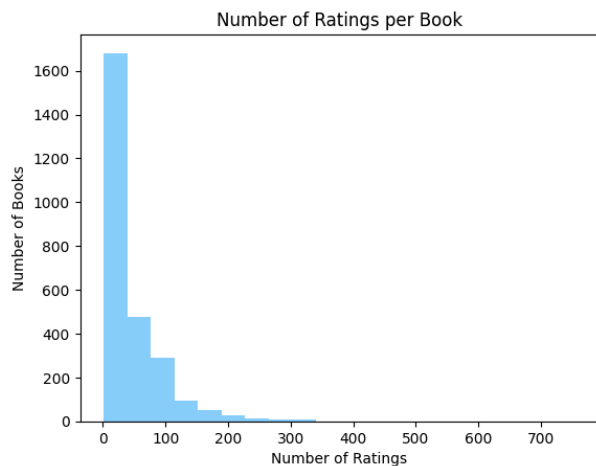
	Year-Of-Publication	User-ID	Book-Rating	Age
count	103943.00	103943.00	103943.00	103943.00
mean	1997.19	139179.81	7.98	35.44
std	5.59	80710.76	0.97	9.81
min	1920.00	243.00	1.00	7.00
25%	1995.00	69042.00	8.00	30.00
50%	1999.00	138189.00	8.00	32.00
75%	2001.00	210792.00	8.00	40.00
max	2010.00	278843.00	10.00	100.00

```
1 filtered_df.head()
```

	ISBN	Book-Title	Book-Author	Year-Of-Publication	Publisher	User-ID	Book-Rating	Location	Age
19	0786868716	The Five People You Meet in Heaven	Mitch Albom	2003	Hyperion	11400	9	ottawa, ontario, canada	49
20	0151008116	Life of Pi	Yann Martel	2002	Harcourt	11400	6	ottawa, ontario, canada	49
21	0671021001	She's Come Undone (Oprah's Book Club)	Wally Lamb	1998	Pocket	11400	8	ottawa, ontario, canada	49

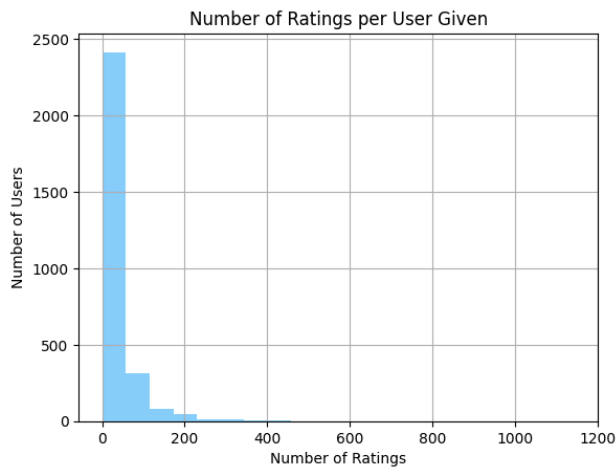
Number of Ratings per Book

```
1 #@title Number of Ratings per Book
2 count_book_ratings = filtered_df['ISBN'].value_counts()
3 plt.hist(count_book_ratings, bins=20, color='lightskyblue')
4 plt.title('Number of Ratings per Book')
5 plt.xlabel('Number of Ratings')
6 plt.ylabel('Number of Books')
7 plt.show()
```



Number of Ratings per User Given

```
1 #@title Number of Ratings per User Given
2 count_user_ratings = filtered_df['User-ID'].value_counts()
3 count_user_ratings.hist(bins=20, color='lightskyblue')
4 plt.title('Number of Ratings per User Given')
5 plt.xlabel('Number of Ratings')
6 plt.ylabel('Number of Users')
7 plt.show()
```

1. Book Recommendation - Collaborative Filtering

Collaborative filtering is a popular technique in recommendation systems. It relies on the assumption that users who have agreed in the past will agree in the future about their preferences.

This method constructs a matrix, with books represented by rows and users by columns, while the matrix entries correspond to the ratings that users have given to books. Given the nature of the dataset, most entries in the matrix are missing and we fill the missing values with 0.

```
1 # Create a pivot table that organizes book ratings with books as rows and users as columns.
2 book_user_rating_pt = filtered_df.pivot_table(index='Book-Title', columns='User-ID', values='Book-Rating')
3 # Fill missing values with 0
4 book_user_rating_pt.fillna(0, inplace=True)
5 book_user_rating_pt.head(5)
```

User-ID	243	254	507	638	643	741	882	929	1211	1424	...	277928	277965	278026	278137	278144	278188	278418	278419
Book-Title																			
1984	0.00	9.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	...	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
1st to Die: A Novel	0.00	0.00	8.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	...	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
24 Hours	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	...	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
2nd	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	...	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

1.1 Cosine Similarity Method

Cosine similarity provides a measure of similarity between two non-zero vectors of an inner product space, it reflects the cosine of the angle between vectors, with a value 1 indicating perfect similarity and a value 0 indicating no similarity. In this method we calculate the cosine similarity between books based on the user ratings, which helps us understand which books are most similar.

```
1 # Compute cosine similarity between books based on user ratings.
2 # Cosine similarity measures how similar the books are in terms of user ratings.
3 similarity_scores = cosine_similarity(book_user_rating_pt)
```

```
1 def cosine_similarity_recommendation(book_title):
2     if book_title not in filtered_df['Book-Title'].values:
3         print("Book title not found in the dataset.")
4         return
5
6     idx = np.where(book_user_rating_pt.index==book_title)[0][0]
7     # Get the list of books based on similarity scores
8     similar_books = sorted(list(enumerate(similarity_scores[idx])), key=lambda x:x[1], reverse=True)[1:10]
9
10    res = []
11    for idx, score in similar_books:
12        book = []
13        book_df = books[books['Book-Title'] == book_user_rating_pt.index[idx]].drop_duplicates('Book-Title')
14        book_details = book_df[['Book-Title', 'Book-Author', 'Year-Of-Publication', 'Publisher']].values.flatten().tolist()
15        res.append(book_details)
16
17    return res
```

```
1 cosine_similarity_recommendation('24 Hours')
```

```

[['The Switch', 'Sandra Brown', 2001, 'Warner Vision'],
 ['Move to Strike', 'Perri O'Shaughnessy', 2001, 'Island'],
 ['Moment of Truth', 'Lisa Scottoline', 2001, 'HarperTorch'],
 ['The Alienist', 'Caleb Carr', 1995, 'Bantam Books'],
 ['The Sigma Protocol', 'Robert Ludlum', 2002, 'St. Martin's Paperbacks'],
 ['Mystic River', 'Dennis Lehane', 2002, 'HarperTorch'],
 ['B Is for Burglar (Kinsey Millhone Mysteries (Paperback))',
  'Sue Grafton',
  1986,
  'Bantam'],
 ['The Simple Truth', 'David Baldacci', 1999, 'Warner Books'],
 ['The Blue Nowhere : A Novel', 'Jeffery Deaver', 2002, 'Pocket']]

```

```
1 consine_similarity_recommendation('1984')
```

```

[['Brave New World', 'Aldous Huxley', 1989, 'Harpercollins'],
 ['Animal Farm', 'George Orwell', 2004, 'Signet'],
 ['Lord of the Flies', 'William Gerald Golding', 1959, 'Perigee Trade'],
 ['The Catcher in the Rye', 'J.D. Salinger', 1991, 'Little, Brown'],
 ['Fahrenheit 451', 'Ray Bradbury', 1994, 'Distribooks Inc'],
 ['Word Freak: Heartbreak, Triumph, Genius, and Obsession in the World of Competitive Scrabble Players',
  'Stefan Fatsis',
  2002,
  'Penguin Books'],
 ['Me Talk Pretty One Day', 'David Sedaris', 2001, 'Back Bay Books'],
 ['To Kill a Mockingbird', 'Harper Lee', 1988, 'Little Brown & Company'],
 ['The Hitchhiker's Guide to the Galaxy', 'Douglas Adams', 1982, 'Pocket']]

```

✓ 1.2 KNN Method

The K-Nearest Neighbors (KNN) is an algorithm used to identify items with the most similarity to a query item. In the context of our Book Recommendation System, it is used to find the k books that have the closest user rating patterns to the given book. The neighbors are selected based on their proximity to the given book using the distance metric defined by the algorithm.

```

1 # Use NearestNeighbors for finding k-nearest items
2 knn_model = NearestNeighbors(metric='cosine', algorithm='brute', n_neighbors=5, n_jobs=-1)
3 knn_model.fit(book_user_rating_pt.values)

```

```

▼ NearestNeighbors
NearestNeighbors(algorithm='brute', metric='cosine', n_jobs=-1)

```

```

1 def knn_recommendation(book_title, model, data, n_recommendations):
2     if book_title not in data.index:
3         print("Book title not found in the dataset.")
4         return
5
6     book_idx = data.index.get_loc(book_title)
7     distances, indices = model.kneighbors(data.iloc[book_idx, :].values.reshape(1, -1), n_neighbors=n_recommendations + 1)
8
9     res = []
10    for i in range(1, len(distances.flatten())):
11        book_info = filtered_df[filtered_df['Book-Title'] == data.index[indices.flatten()[i]]].drop_duplicates('Book-Title')
12        book_details = book_info[['Book-Title', 'Book-Author', 'Year-Of-Publication', 'Publisher']].values.flatten().tolist()
13        res.append(book_details)
14
15    return res

```

```
1 knn_recommendation('24 Hours', knn_model, book_user_rating_pt, 10)
```

```

[['The Switch', 'Sandra Brown', 2001, 'Warner Vision'],
 ['Move to Strike', 'Perri O'Shaughnessy', 2001, 'Island'],
 ['Moment of Truth', 'Lisa Scottoline', 2001, 'HarperTorch'],
 ['The Alienist', 'Caleb Carr', 1995, 'Bantam Books'],
 ['The Sigma Protocol', 'Robert Ludlum', 2002, 'St. Martin's Paperbacks'],
 ['Mystic River', 'Dennis Lehane', 2002, 'HarperTorch'],
 ['B Is for Burglar (Kinsey Millhone Mysteries (Paperback))',
  'Sue Grafton',
  1986,
  'Bantam'],
 ['The Simple Truth', 'David Baldacci', 1999, 'Warner Books'],
 ['The Blue Nowhere : A Novel', 'Jeffery Deaver', 2002, 'Pocket'],
 ['A Map of the World', 'Jane Hamilton', 1999, 'Anchor Books/Doubleday']]

```

```
1 knn_recommendation('1984', knn_model, book_user_rating_pt, 10)
```

```

[['Brave New World', 'Aldous Huxley', 1989, 'Harpercollins'],
 ['Animal Farm', 'George Orwell', 2004, 'Signet'],
 ['Lord of the Flies', 'William Gerald Golding', 1959, 'Perigee Trade'],
 ['The Catcher in the Rye', 'J.D. Salinger', 1991, 'Little, Brown'],
 ['Fahrenheit 451', 'Ray Bradbury', 1994, 'Distribooks Inc'],
 ['Word Freak: Heartbreak, Triumph, Genius, and Obsession in the World of Competitive Scrabble Players',
  'Stefan Fatsis',
  2002,
  'Penguin Books'],
 ['Me Talk Pretty One Day', 'David Sedaris', 2001, 'Back Bay Books'],
 ['To Kill a Mockingbird', 'Harper Lee', 1988, 'Little Brown & Company'],
 ['The Hitchhiker's Guide to the Galaxy', 'Douglas Adams', 1982, 'Pocket'],
 ['Fast Food Nation: The Dark Side of the All-American Meal',
  'Eric Schlosser',
  2002,
  'Perennial']]

```

2. Book Recommendation - Clustering

The clustering algorithm categorizes users into segments based on their book ratings. This approach treats the recommendation as a classification problem. Users in the same cluster share similar preferences and behaviors, and are more alike to each other than in other clusters. By utilizing the clustering algorithm, the system can recommend books that resonate with the specific interests of each user group.

```
1 # Create a pivot table that organizes book ratings with users as rows and books as columns.
2 user_book_rating_pt = filtered_df.pivot_table(index='User-ID', columns='Book-Title', values='Book-Rating')
3
4 # Fill missing values with 0
5 user_book_rating_pt.fillna(0,inplace=True)

1 # Perform clustering and assign users to clusters
2 kmeans = KMeans(n_clusters=100, random_state=42).fit(user_book_rating_pt)
3 user_book_rating_pt['Cluster'] = kmeans.labels_

1 user_book_rating_pt.head(5)
```

Book-Title	1984	1st to Die: A Novel	24 Hours	2nd Chance	Blondes	4 Mathematical Genius and Nobel Laureate John Nash	A Beautiful Mind: The Life of	A Bend in the Road	A Case of Need	A Child Called \It\": One Child's Courage to Survive"	A Civil ...	Wizard and Glass (The Dark Tower, Book 4)	Women Who Run with the Wolves	Word Freak: Heartbreak, Triumph, Genius, and Obsession in the World of Competitive Scrabble Players	Wu
User-ID															
243	0.00	0.00	0.00	8.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
254	9.00	0.00	0.00	0.00	0.00	0.00	0.00	8.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
507	0.00	8.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
638	0.00	0.00	0.00	9.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
643	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

5 rows x 894 columns

```
1 def cluster_recommendation(user_id, num_recommendations=5):
2     if user_id not in user_book_rating_pt.index:
3         return f"User ID {user_id} not found."
4
5     # Find the cluster of the user
6     cluster = user_book_rating_pt.loc[user_id, 'Cluster']
7     # Get all users in the cluster
8     cluster_users = user_book_rating_pt[user_book_rating_pt['Cluster'] == cluster]
9
10    # Get average book ratings of that cluster
11    book_ratings = cluster_users.replace(0, pd.NA).mean().sort_values(ascending=False)
12
13    # Filter out books the given user has already rated
14    booksRated = user_book_rating_pt.columns[user_book_rating_pt.loc[user_id] > 0]
15    recommendations = book_ratings.drop(labels=booksRated).head(num_recommendations)
16
17    print("Books enjoyed by users with similar tastes:")
18    return recommendations
```

```
1 cluster_recommendation(2276, 10)
```

Books enjoyed by users with similar tastes:

Book-Title	
The Little Prince	10.00
Chicken Soup for the Teenage Soul (Chicken Soup for the Soul)	10.00
Siddhartha	9.50
The Last Time They Met : A Novel	9.00
Forever... : A Novel of Good and Evil, Love and Hope	9.00
Atlas Shrugged	9.00
Catering to Nobody	9.00
Naked	9.00
Sisterhood of the Traveling Pants	9.00
Lucky Man: A Memoir	9.00
dtype: object	

```
1 cluster_recommendation(3363, 10)
```

Books enjoyed by users with similar tastes:

Book-Title	
Dark Rivers of the Heart	10.00
Slaughterhouse Five or the Children's Crusade: A Duty Dance With Death	9.00
The Sparrow	9.00
Mind Prey	9.00
Dragon Tears	9.00

Where the Red Fern Grows	8.89
Charlotte's Web (Trophy Newbery)	8.83
Interpreter of Maladies	8.77
Ender's Game (Ender Wiggins Saga (Paperback))	8.69
Anne Frank: The Diary of a Young Girl	8.69
dtype: object	

3. Book Recommendation - Content-Based Using Word2Vec

This approach uses the Word2Vec model to perform content-based filtering. Word2Vec is utilized to understand and quantify the semantic relationships between words within book titles and their associated genres.

By training the Word2Vec model on the book titles text data, this approach generates word embeddings that capture the essence of each book's thematic content, and predicts the genre of each book.

This recommendation approach identifies the top 10 books with the highest average ratings within the same genre as the given book.

```

1 def clean_text(text):
2     text = re.sub(r'^a-zA-Z\s', '', text).lower()
3     return text

1 # Define genres
2 genres = ['Romance', 'Mystery', 'Science', 'History', 'Fiction', 'Horror', 'Thriller', 'Other']
3
4 # Split the title words, prepare words for training
5 words = [[genre.lower() for genre in genres]]
6 for title in filtered_df['Book-Title']:
7     words.append(clean_text(title).split())
8
9 # Train the Word2Vec model
10 word2vec_model = Word2Vec(words, vector_size=100, window=5, min_count=1, sg=1)

1 def get_most_similar(title_vector):
2     # Initialize max similarity and result
3     max_similarity = float('-inf')
4     res = None
5
6     # Iterate over genres and find the most similar genre
7     for genre in genres:
8         genre_vector = word2vec_model.wv[genre.lower()]
9         similarity = cosine_similarity([title_vector], [genre_vector])[0][0]
10        if similarity > max_similarity:
11            res = genre
12            max_similarity = similarity
13    return res
14
15 def predict_genre(title):
16     # Clean up the title and split into words
17     # Filter title words that are present in the Word2Vec model vocabulary
18     words = clean_text(title).split()
19     title_words = []
20     for word in words:
21         if word in word2vec_model.wv:
22             title_words.append(word)
23
24     # If there is no word left, return 'Other' genre
25     if not title_words:
26         return 'Other'
27
28     # Calculate average vector of the title words
29     avg_vector = sum([word2vec_model.wv[w] for w in title_words]) / len(title_words)
30
31     # Find the most similar genre
32     return get_most_similar(avg_vector)

1 predict_genre('A Map of the World')

'History'

1 predict_genre('Horror')

'Horror'

1 # Apply the function and get all books predicted
2 filtered_df['Genre'] = filtered_df['Book-Title'].swifter.apply(predict_genre)
3 filtered_df.head()
```

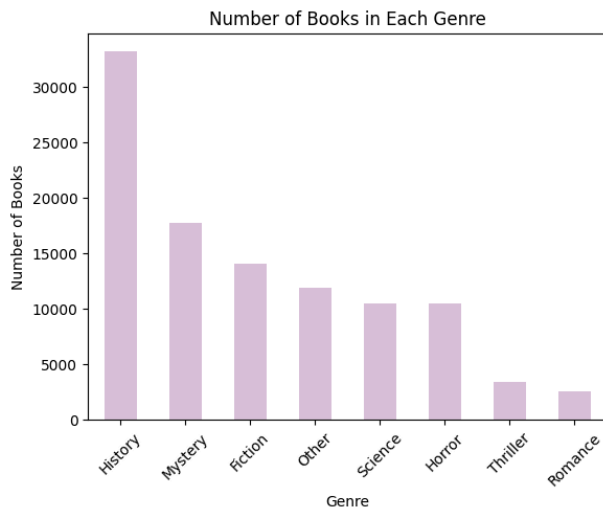
```
<ipython-input-346-53227ec46956>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-swifter对象
 filtered_df['Genre'] = filtered_df['Book-Title'].swifter.apply(predict_genre)

	ISBN	Book-Title	Book-Author	Year-Of-Publication	Publisher	User-ID	Book-Rating	Location	Age	Genre
19	0786868716	The Five People You Meet in Heaven	Mitch Albom	2003	Hyperion	11400	9	ottawa, ontario, canada	49	History

Number of Books in Each Genre

```
1 #@title Number of Books in Each Genre
2 count_genres = filtered_df.groupby('Genre')['Book-Title'].count().sort_values(ascending=False)
3 count_genres.plot(kind='bar', color='thistle')
4 plt.title('Number of Books in Each Genre')
5 plt.xlabel('Genre')
6 plt.ylabel('Number of Books')
7 plt.xticks(rotation=45)
8 plt.show()
```



```
1 def content_based_recommendation(book_title):
2     if book_title not in filtered_df['Book-Title'].values:
3         print("Book title not found in the dataset.")
4         return
5
6     genre = filtered_df[filtered_df['Book-Title'] == book_title]['Genre'].iloc[0]
7     books_same_genre = filtered_df[filtered_df['Genre'] == genre].groupby('Book-Title')['Book-Rating'].mean().sort_values(ascen
8
9     recommend_books = books_same_genre.head(10).index.tolist()
10
11     res = []
12     for title in recommend_books:
13         book_info = filtered_df[(filtered_df['Book-Title'] == title) & (filtered_df['Genre'] == genre)].drop_duplicates('Book-Tit
14         book_details = book_info[['Book-Title', 'Book-Author', 'Year-Of-Publication', 'Publisher']].values.flatten().tolist()
15         res.append(book_details)
16
17     print(f"More books in genre '{genre}':")
18     return res
```

```
1 content_based_recommendation('Animal Farm')
```

```
More books in genre 'Mystery':
[["Ender's Game (Ender Wiggins Saga (Paperback))",
  'Orson Scott Card',
  1994,
  'Tor Books'],
 ['Anne of Avonlea (Anne of Green Gables Novels (Paperback))',
  'L.M. MONTGOMERY',
  1984,
  'Bantam Classics'],
```