

In [10]:

```
1 import numpy as np
2 import pandas as pd
3
4 import warnings
5 warnings.filterwarnings("ignore")
6
7 books = pd.read_csv("BX-Books.csv", sep=';', encoding="latin-1", on_bad_lines='skip')
8 users = pd.read_csv("BX-Users.csv", sep=';', encoding="latin-1", on_bad_lines='skip')
9 ratings = pd.read_csv("BX-Book-Ratings.csv", sep=';', encoding="latin-1", on_bad_lines='skip')
10
11 books.head()
```

Out[10]:

	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.com/ima
1	0002005018	Clara Callan	Richard Bruce Wright	2001	HarperFlamingo Canada	http://images.amazon.com/images/P/0002005018.0...	http://images.amazon.com/ima
2	0060973129	Decision in Normandy	Carlo D'Este	1991	HarperPerennial	http://images.amazon.com/images/P/0060973129.0...	http://images.amazon.com/ima
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.com/ima
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.com/ima

```
In [ ]: 1 # Preprocessing Data
        2 # Now in the books file, we have some extra columns which are not required for our task like image L
        3 # And we will rename the columns of each file as the name of the column contains space,
        4 # and uppercase letters so we will correct as to make it easy to use.
```

```
In [11]: 1 books = books[['ISBN', 'Book-Title', 'Book-Author', 'Year-Of-Publication', 'Publisher']]
        2 books.rename(columns = {'Book-Title':'title', 'Book-Author':'author', 'Year-Of-Publication':'year',
        3 users.rename(columns = {'User-ID':'user_id', 'Location':'location', 'Age':'age'}, inplace=True)
        4 ratings.rename(columns = {'User-ID':'user_id', 'Book-Rating':'rating'}, inplace=True)
```

```
In [ ]: 1 # when you run the above code we can see only 105283 peoples have given a rating among 278000.
        2 # Now we will extract the user ids who have given more than 200 ratings and when we will have
        3 # user ids we will extract the ratings of only this user id from the rating dataframe.
```

```
In [12]: 1 x = ratings['user_id'].value_counts() > 200
        2 y = x[x].index #user_ids
        3 print(y.shape)
        4 ratings = ratings[ratings['user_id'].isin(y)]
```

(899,)

```
In [ ]: 1 # So 900 users are there who have given 5.2 lakh rating and this we want.
        2 # Now we will merge ratings with books on basis of ISBN so that we will get the
        3 # rating of each user on each book id and the user who has not rated that book id the value will be
```

```
In [13]: 1 rating_with_books = ratings.merge(books, on='ISBN')
          2 rating_with_books.head()
```

Out[13]:

	user_id	ISBN	rating		title	author	year		publisher
0	277427	002542730X	10		Politically Correct Bedtime Stories: Modern Ta...	James Finn Garner	1994		John Wiley & Sons Inc
1	3363	002542730X	0		Politically Correct Bedtime Stories: Modern Ta...	James Finn Garner	1994		John Wiley & Sons Inc
2	11676	002542730X	6		Politically Correct Bedtime Stories: Modern Ta...	James Finn Garner	1994		John Wiley & Sons Inc
3	12538	002542730X	10		Politically Correct Bedtime Stories: Modern Ta...	James Finn Garner	1994		John Wiley & Sons Inc
4	13552	002542730X	0		Politically Correct Bedtime Stories: Modern Ta...	James Finn Garner	1994		John Wiley & Sons Inc

```
In [ ]: 1 # Extract books that have received more than 50 ratings.
          2 # Now dataframe size has decreased and we have 4.8 lakh because when we merge the dataframe,
          3 # all the book id-data we were not having. Now we will count the rating of each book so we will
          4 # group data based on title and aggregate based on rating.
```

```
In [14]: 1 number_rating = rating_with_books.groupby('title')['rating'].count().reset_index()
          2 number_rating.rename(columns= {'rating':'number_of_ratings'}, inplace=True)
          3 final_rating = rating_with_books.merge(number_rating, on='title')
          4 final_rating.shape
          5 final_rating = final_rating[final_rating['number_of_ratings'] >= 50]
          6 final_rating.drop_duplicates(['user_id', 'title'], inplace=True)
```

```
In [ ]: 1 # Create Pivot Table
          2 # As we discussed above we will create a pivot table where columns will be user ids,
          3 # the index will be book title and the value is ratings. And the user id who has not
          4 # rated any book will have value as NAN so impute it with zero.
```

```
In [15]: 1 book_pivot = final_rating.pivot_table(columns='user_id', index='title', values="rating")
          2 book_pivot.fillna(0, inplace=True)
```

```
In [ ]: 1 # Modeling
        2 # We have prepared our dataset for modeling. we will use the nearest neighbors algorithm which is
        3 # the same as K nearest which is used for clustering based on euclidian distance.
        4
        5 # But here in the pivot table, we have lots of zero values and on clustering,
        6 # this computing power will increase to calculate the distance of zero values so we will convert
        7 # the pivot table to the sparse matrix and then feed it to the model.
```

```
In [16]: 1 from scipy.sparse import csr_matrix
        2 book_sparse = csr_matrix(book_pivot)
```

```
In [17]: 1 from sklearn.neighbors import NearestNeighbors
        2 model = NearestNeighbors(algorithm='brute')
        3 model.fit(book_sparse)
```

```
Out[17]: ▼      NearestNeighbors
          NearestNeighbors(algorithm='brute')
```

```
In [18]: 1 distances, suggestions = model.kneighbors(book_pivot.iloc[237, :].values.reshape(1, -1))
        2
```

```
In [19]: 1 for i in range(len(suggestions)):
        2     print(book_pivot.index[suggestions[i]])
```

```
Index(['Harry Potter and the Chamber of Secrets (Book 2)',
      'Harry Potter and the Prisoner of Azkaban (Book 3)',
      'Harry Potter and the Goblet of Fire (Book 4)',
      'Harry Potter and the Sorcerer's Stone (Book 1)', 'Exclusive'],
      dtype='object', name='title')
```

```
In [ ]: 1
```

