

In [13]: `import numpy as np
import pandas as pd
import warnings
warnings.filterwarnings('ignore')`

Read the data using pandas in DataFrame df_user
To map byte values directly to the first 256 Unicode code points, use the "Latin-1" encoding. This is the closest equivalent Python 3 offers to the permissive Python 2 text handling model.

In [14]: `df_user = pd.read_csv('BX-Users.csv',encoding='latin-1')`

In [15]: `df_user.shape`

Out[15]: (278859, 3)

In [16]: `df_user.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 278859 entries, 0 to 278858
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   user_id     278859 non-null  object
1   Location    278858 non-null  object
2   Age         168896 non-null  float64
dtypes: float64(1), object(2)
memory usage: 6.4+ MB
```

In [17]: `# Checking for Null Values.
df_user.isnull().sum()`

Out[17]: `user_id 0
Location 1
Age 119763
dtype: int64`

In [18]: `# Dropping the Null Values.
df_user1=df_user.dropna()`

In [19]: `df_user1.isnull().sum()`

Out[19]: `user_id 0
Location 0
Age 0
dtype: int64`

Read the books Data and explore

In [20]: `df_books = pd.read_csv('BX-Books.csv', encoding='latin-1')`

In [21]: `df_books.shape`

Out[21]: (271379, 5)

Reading the data where ratings are given
We will read only first 10000 rows otherwise, Out Of Memory error can occur.

In [22]: `df_ratings = pd.read_csv('BX-Book-Ratings.csv',encoding='latin-1',nrows=10000)`

Using 'describe()' function
It is used to view some basic statistical details like percentile, mean, std.

In [23]: `df_ratings.describe()`

Out[23]:

	user_id	rating
count	10000.000000	10000.000000
mean	265844.379600	1.974700
std	56937.189618	3.424884
min	2.000000	0.000000
25%	277478.000000	0.000000
50%	278418.000000	0.000000
75%	278418.000000	4.000000
max	278854.000000	10.000000

Merge the dataframes
For all practical purposes, User Master Data is not required. So, ignore dataframe df_user

In [24]: `df_final = pd.merge(df_ratings,df_books,on='isbn')`

In [25]: `df_final.head()`

Out[25]:

	user_id	isbn	rating	book_title	book_author	year_of_publication	publisher
0	276725	034545104X	0	Flesh Tones: A Novel	M. J. Rose	2002	Ballantine Books
1	276726	155061224	5	Rites of Passage	Judith Rae	2001	Heinle
2	276727	446520802	0	The Notebook	Nicholas Sparks	1996	Warner Books
3	278418	446520802	0	The Notebook	Nicholas Sparks	1996	Warner Books
4	276729	052165615X	3	Help!: Level 1	Philip Prowse	1999	Cambridge University Press

Checking for unique users and books
Here we are using 'nunique()' function that returns the Series with the number of distinct observations over the requested axis.

In [26]: `# Code For checking number of unique users and books.
n_users = df_final.user_id.nunique()
n_books = df_final.isbn.nunique()

print('Num. of Users: ' + str(n_users))
print('Num of Books: ' +str(n_books))

Num. of Users: 828
Num of Books: 8051`

Convert ISBN variable to numeric type in order

In [27]: `# Convert and print length of isbn list
isbn_list = df_final.isbn.unique()
print(" Length of isbn List:", len(isbn_list))
def get_isbn_numeric_id(isbn):
 #print (" isbn is:" , isbn)
 itemindex = np.where(isbn_list==isbn)
 return itemindex[0][0]`

Length of isbn List: 8051

Convert user_id variable to numeric type in order
This is formatted as code.

In [28]: `# Convert and print length of user_id list
userid_list = df_final.user_id.unique()
print(" Length of user_id List:", len(userid_list))
def get_user_id_numeric_id(user_id):
 #print (" isbn is:" , isbn)
 itemindex = np.where(userid_list==user_id)
 return itemindex[0][0]`

Length of user_id List: 828

Convert both user_id and isbn to ordered list i.e. from 0..n-1

In [29]: `df_final['user_id_order'] = df_final['user_id'].apply(get_user_id_numeric_id)`

In [30]: `df_final['isbn_id'] = df_final['isbn'].apply(get_isbn_numeric_id)
df_final.head()`

Out[30]:

	user_id	isbn	rating	book_title	book_author	year_of_publication	publisher	user_id_order	isbn_id
0	276725	034545104X	0	Flesh Tones: A Novel	M. J. Rose	2002	Ballantine Books	0	0
1	276726	155061224	5	Rites of Passage	Judith Rae	2001	Heinle	1	1
2	276727	446520802	0	The Notebook	Nicholas Sparks	1996	Warner Books	2	2
3	278418	446520802	0	The Notebook	Nicholas Sparks	1996	Warner Books	3	2
4	276729	052165615X	3	Help!: Level 1	Philip Prowse	1999	Cambridge University Press	4	3

Re-index columns to build matrix

In [31]: `# Reindexing the columns
new_col_order = ['user_id_order', 'isbn_id', 'rating', 'book_title', 'book_author','year_of_publication','publisher','isbn','user_id']
df_final = df_final.reindex(columns=new_col_order)
df_final.head()`

Out[31]:

	user_id_order	isbn_id	rating	book_title	book_author	year_of_publication	publisher	isbn	user_id
0	0	0	0	Flesh Tones: A Novel	M. J. Rose	2002	Ballantine Books	034545104X	276725
1	1	1	5	Rites of Passage	Judith Rae	2001	Heinle	155061224	276726
2	2	2	0	The Notebook	Nicholas Sparks	1996	Warner Books	446520802	276727
3	3	2	0	The Notebook	Nicholas Sparks	1996	Warner Books	446520802	278418
4	4	3	3	Help!: Level 1	Philip Prowse	1999	Cambridge University Press	052165615X	276729

Train Test Split

Recommendation Systems are difficult to evaluate, but we will still learn how to evaluate them. In order to do this, will split our data into two sets. However, we won't do our classic X_train,X_test,y_train,y_test split. Instead, we can actually just segment the data into two sets of data:

Importing train_test_split model

In [32]: `# Importing train_test_split model for splitting the data into train and test set.
from sklearn.model_selection import train_test_split
train_data, test_data = train_test_split(df_final, test_size=0.20)`

Approach: We Will Use Memory-Based Collaborative Filtering

Memory-Based Collaborative Filtering approaches can be divided into two main sections: **user-item filtering** and **item-item filtering**.

A *user-item filtering* will take a particular user, find users that are similar to that user based on similarity of ratings, and recommend items that those similar users liked.

In contrast, *item-item filtering* will take an item, find users who liked that item, and find other items that those users or similar users also liked. It takes items as input and outputs other items as recommendations.

- *Item-Item Collaborative Filtering*: "Users who liked this item also liked ..."
- *User-Item Collaborative Filtering*: "Users who are similar to you also liked ..."

In both cases, we will create a user-book matrix which is built from the entire dataset.

Since we have split the data into testing and training, we will need to create two [828 x 8051] matrices (all users by all books). This is going to be a very large matrix.

The training matrix contains 80% of the ratings and the testing matrix contains 20% of the ratings.

Create two user-book matrix for training and testing

 Indented block

In [33]: `# Create user-book matrix for training
train_data_matrix = np.zeros((n_users, n_books))
for line in train_data.iteruples():
 train_data_matrix[line[1]-1, line[2]-1] = line[3]

Create user-book matrix for testing
test_data_matrix = np.zeros((n_users, n_books))
for line in test_data.iteruples():
 test_data_matrix[line[1]-1, line[2]-1] = line[3]`

Import Pairwise Model
we can use the *pairwise_distances* function from *sklearn* to calculate the cosine similarity. Note, the output will range from 0 to 1 since the ratings are all positive.

In [34]: `# Importing pairwise_distances function
from sklearn.metrics.pairwise import pairwise_distances
user_similarity = pairwise_distances(train_data_matrix, metric='cosine')
item_similarity = pairwise_distances(train_data_matrix.T, metric='cosine')`

In [35]: `user_similarity`

Out[35]: `array([[0., 1., 1., ..., 1., 1., 1.],
 [1., 0., 1., ..., 1., 1., 1.],
 [1., 1., 0., ..., 1., 1., 1.],
 ...,
 [1., 1., 1., ..., 0., 1., 1.],
 [1., 1., 1., ..., 1., 0., 1.],
 [1., 1., 1., ..., 1., 1., 0.]])`

Make predictions

In [36]: `# Defining custom function to make predictions
def predict(ratings, similarity, type='user'):
 if type == 'user':
 mean_user_rating = ratings.mean(axis=1)
 # We will use np.newaxis so that mean_user_rating has same format as ratings.
 ratings_diff = (ratings - mean_user_rating[:, np.newaxis])
 pred = mean_user_rating[:, np.newaxis] + similarity.dot(ratings_diff) / np.array([np.abs(similarity).sum(axis=1)]).T
 elif type == 'item':
 pred = ratings.dot(similarity) / np.array([np.abs(similarity).sum(axis=1)])
 return pred`

In [37]: `item_prediction = predict(train_data_matrix, item_similarity, type='item')
user_prediction = predict(train_data_matrix, user_similarity, type='user')`

In [38]: `print(item_prediction)`

```
[[0. 0.00062112 0.0008212 ... 0.00082167 0.00062112 0.00062112]
 [0. 0. ... 0. 0. ]
 [0.0642236 0.0642236 0.06423158 ... 0.06428079 0.0642236 0.0642236 ]
 ...
 [0. 0. 0. ... 0. 0. 0. ]
 [0. 0. 0. ... 0. 0. 0. ]
 [0. 0. 0. ... 0. 0. 0. ]]
```

In [39]: `print(user_prediction)`

```
[[-0.00140624 -0.00140624 0.00222133 ... 0.00947647 -0.00140624
 -0.00140624]
 [ 0.00401792 -0.00202803 0.00159954 ... 0.00885468 -0.00202803
 -0.00202803]
 [ 0.0083137 0.06226614 0.06589468 ... 0.07315176 0.06226614
 0.06226614]
 ...
 [ 0.00401792 -0.00202803 0.00159954 ... 0.00885468 -0.00202803
 -0.00202803]
 [ 0.00401792 -0.00202803 0.00159954 ... 0.00885468 -0.00202803
 -0.00202803]
 [ 0.00401792 -0.00202803 0.00159954 ... 0.00885468 -0.00202803
 -0.00202803]]
```

Evaluation

There are many evaluation metrics, but one of the most popular metric used to evaluate accuracy of predicted ratings is *Root Mean Squared Error (RMSE)*.

Since, we only want to consider predicted ratings that are in the test dataset, we will filter out all other elements in the prediction matrix with: `prediction[ground_truth.nonzero()]`.

In [40]: `# Importing RMSE function
from sklearn.metrics import mean_squared_error
from math import sqrt

Defining custom function to filter out elements with ground_truth.nonzero
def rmse(prediction, ground_truth):
 prediction = prediction[ground_truth.nonzero()].flatten()
 ground_truth = ground_truth[ground_truth.nonzero()].flatten()
 return sqrt(mean_squared_error(prediction, ground_truth))`

Printing RMSE value for user based and item based collaborative filtering

In [41]: `print('User-based CF RMSE: ' + str(rmse(user_prediction, test_data_matrix)))
print('Item-based CF RMSE: ' + str(rmse(item_prediction, test_data_matrix)))`

