# Long term project

컴퓨터소프트웨어학과 2017029470 김종원

## 1. Environment

- Mac OS (Monterey) M1 chip
- Python 3.10.3 (release March 16, 2022)

## 2. How to compile

Before compiling dt.py, python version 3 must be installed in your system.

```
wonnx@wonnx recommender_system_project % python3 recommender.py data-1/u1.base data-1/u1.test
wonnx@wonnx recommender_system_project % python3 recommender.py data-1/u2.base data-1/u2.test
wonnx@wonnx recommender_system_project % python3 recommender.py data-1/u3.base data-1/u3.test
wonnx@wonnx recommender_system_project % python3 recommender.py data-1/u4.base data-1/u4.test
wonnx@wonnx recommender_system_project % python3 recommender.py data-1/u5.base data-1/u5.test
```

- Execution file name: recommender.py
- Training data file name: u#.test.txt
- Test data file name: u#.test.txt
- Output file name: u#.base_prediction.txt

## 3. Summary of algorithm

I used the collaborate filtering method to predict the ratings of movies in test data. Collaborative filtering refers to a technology that predicts itself based on taste information collected from many users. Based on the user-user similarity neighbors that most similar to the corresponding user were extracted, and ratings were predicted based on the rates of the extracted neighbors.

- Using the pandas library, the ratings of movies left by users are transformed into a rating matrix.
- Based on the rating matrix, the Pearson correlation coefficient between user-user is obtained and converted into a similarity matrix.
- After sorting the similarity matrix, the neighbors corresponding to the top 5% of all users are obtained, and then the rating is predicted through a weighted sum.

## 4. Detailed description of codes

### A. get_rating_matrix

```
5    # Transforms the given data into a rating matrix form
6    # using the pivot_table of the pandas library.
7    def get_rating_matrix():
8        matrix = np.loadtxt(sys.argv[1], dtype='int', usecols=(0, 1, 2))
9        matrix = pd.DataFrame({'user': matrix[:, 0],
10                               'item': matrix[:, 1],
11                               'rating': matrix[:, 2]})
12       rating_matrix = matrix.pivot_table('rating', index='user',columns='item', fill_value=0)
13       return rating_matrix
```

This is a function that transforms the given training data into a rating matrix. In the training data, the column corresponding to the time stamp was not used.

- Convert training data into dataFrame format of pandas library.
- Convert the dataFrame named **matrix** into the pivot_table type **rating_matrix**.

### B. get_similarity_matrix

```
15   # Create a similarity matrix between users
16   # using Pearson correlation coefficient (PCC) to find neighbors.
17   def get_similarity_matrix(rating_matrix):
18       similarity_matrix = (rating_matrix.T).corr(method='pearson')
19       return similarity_matrix
```

Based on the previously obtained rating matrix, it is a function that calculates the similarity between users and users though the Pearson correlation coefficient and stored as a similarity matrix.

- Using the built-in correlation method of pandas dataFrame, calculate similarity between users.

### C. predict

```
40   def predict(rating_matrix, similarity_matrix):
41       # Get average rate and neighbors of the each user
42       average = dict()
43       neighbors = dict()
44       for uid in list(rating_matrix.index):
45           rating = list(rating_matrix.loc[uid])
46           cnt = len(list(filter(lambda x: x!=0, rating)))
47           if cnt != 0: average[uid] = sum(rating)/cnt
48           else: average[uid] = 0
49           neighbor = similarity_matrix[uid].sort_values(ascending=False)
50           neighbor = neighbor.reset_index().values.tolist()
51           neighbors[uid] = neighbor
52
53       output = sys.argv[1] + "_prediction.txt"
54       with open(sys.argv[2], 'r') as test_file, open(output, 'w') as file:
55           while True:
56               line = test_file.readline().strip()
57               if not line or len(line)==0: break
58               line = line.split('\t')
59               prediction = estimate(int(line[0]), int(line[1]), rating_matrix, average, neighbors)
60               result = line[0] + '\t' + line[1] + '\t' + str(prediction) + '\n'
61               file.write(result)
```

Average ratings of movies evaluated by each user and neighbors of each user are obtained. Using the average, neighbor, and rating matrix, we predict the user's rating for a movie that has not yet been rated. In addition, the output file name is set differently according to the input file, and the predicted rating is

exported according to the format.

- **Average** and **neighbors** are in the form of dictionary, and each key contains average rating of the movies for the user and PCC for all users.

## D. estimate

```python
21    # Estimate the rating of the test data
22    # By user-based collaborate filtering
23    def estimate(user, item, rating_matrix, average, neighbors):
24        if item not in list(rating_matrix.columns): return 1.0
25        weight = weighted_sum = count = 0
26        for neighbor in neighbors[user]:
27            rate = rating_matrix.loc[neighbor[0]][item]
28            if rate == 0: continue
29            if(count>len(neighbors[user])/20): break
30            weighted_sum += neighbor[1] * (rate - average[user])
31            weight += neighbor[1]
32            count += 1
33        if weight == 0: return average[user]
34        prediction = average[user] + weighted_sum / weight
35        if prediction > 5: prediction = 5.0
36        elif prediction < 1: prediction = 1.0
37        return prediction
```

5% of all users are set as neighbors, and if the neighbor has previously rated the movie, it is multiplied by similarity and weighted summed.

- If the movie does not exist in the existing rating matrix, the rating is predicted using the average rating value of the user.
- If the prediction result is less than 1, it is corrected to 1, and if it is greater than 5, corrected to 5.

## 5. Testing result

```
C:\Users\labor\test>PA4.exe u1
the number of ratings that didn't be predicted: 0
the number of ratings that were unproperly predicted [ex. >=10, <0, NaN, or format errors]: 0
If the counted number is large, please check your codes again.

The bigger value means that the ratings are predicted more incorrectly
RMSE: 1.019897
```

```
C:\Users\labor\test>PA4.exe u2
the number of ratings that didn't be predicted: 0
the number of ratings that were unproperly predicted [ex. >=10, <0, NaN, or format errors]: 0
If the counted number is large, please check your codes again.

The bigger value means that the ratings are predicted more incorrectly
RMSE: 1.013268
```

```
C:\Users\labor\test>PA4.exe u3
the number of ratings that didn't be predicted: 0
the number of ratings that were unproperly predicted [ex. >=10, <0, NaN, or format errors]: 0
If the counted number is large, please check your codes again.

The bigger value means that the ratings are predicted more incorrectly
RMSE: 1.005027
```

```
C:\Users\labor\test>PA4.exe u4
the number of ratings that didn't be predicted: 0
the number of ratings that were unproperly predicted [ex. >=10, <0, NaN, or format errors]: 0
If the counted number is large, please check your codes again.

The bigger value means that the ratings are predicted more incorrectly
RMSE: 1.00327
```

```
C:\Users\labor\test>PA4.exe u5
the number of ratings that didn't be predicted: 0
the number of ratings that were unproperly predicted [ex. >=10, <0, NaN, or format errors]: 0
If the counted number is large, please check your codes again.

The bigger value means that the ratings are predicted more incorrectly
RMSE: 1.009369
```

For testing, I used a measure called RMSE (Root Mean Square Error). The bigger value means that the rating are predicted more incorrectly. All five testing results came out close to 1, indicating that the user actually evaluated the movie and the predicted value differed by about 1.