

Web Science: Assignment #6

Alexander Nwala

Mohd. Nauman Siddique

Sunday, March 31, 2019

Contents

About the data set used in the assignment	3
Problem 1	4
Problem 2	7
Problem 3	9
Problem 4	9

3—23—M—writer—32067

4—24—M—technician—43537

5—33—F—other—15213

The code for reading from the u.data and u.item files and creating recommendations is described in the book Programming Collective Intelligence. Feel free to modify the PCI code to answer the following questions.

Problem 1

Find 3 users who are closest to you in terms of age, gender, and occupation. For each of those 3 users:

1. what are their top 3 favorite films?
2. bottom 3 least favorite films?

Based on the movie values in those 6 tables (3 users X (favorite + least)), choose a user that you feel is most like you. Feel free to note any outliers (e.g., "I mostly identify with user 123, except I did not like "Ghost" at all").

This user is the "substitute you".

SOLUTION

The solution to the problem is as below:

1. **Loading Movie Lens Dataset:** I read the u.data, u.user and u.item table from the Movie Lens data set and loaded them into separate json files.
2. **Find closest users :** We find the top three users that manage my category. We then find the highest rated and least rated movies for each of the three users.
3. **Most matching user:** The user 774 matches my choice of movies from all the three users. Figure 1 shows the results for similar users and their least favourite and most favourite movies.

```
import json
import csv

5 def read_users_table():
    path_data_set = "/home/msiddique/WSDL_Work/WebScience/Assignment6/ml-100k/"
    file_open = open(path_data_set + "u.user", "r")
    list_json = []
    reader = csv.reader(file_open, delimiter="|")
10    for row in reader:
        insert_row = {'user': row[0], 'age': row[1], 'gender': row[2], 'occupation':
            row[3], 'zip': row[4]}
        list_json.append(insert_row)
    file_json = open("Users.json", "w")
    json.dump(list_json, file_json)
15    file_open.close()
    file_json.close()
```

```
def read_items_table():
20 path_data_set = "/home/msiddique/WSDL_Work/WebScience/Assignment6/ml-100k/"
    # ratings = pd.read_table(path_dataset + 'u.data', sep=':', names=['user', '
        movie', 'rating', 'time'])
    file_open = open(path_data_set + "u.item", "r", encoding = "ISO-8859-1")
    list_json = []
    reader = csv.reader(file_open, delimiter="|")
25 for row in reader:
    print(row)
    insert_row = {"movie id": row[0], "movie title": row[1], "release date":
        row[2], "video release date": row[3],
            "IMDb URL": row[4], "unknown": row[5], "Action": row[6], "
                Adventure": row[7], "Animation": row[8],
                    "Children's": row[9], "Comedy": row[10], "Crime": row[11], "
                        Documentary": row[12],
30 "Drama": row[13], "Fantasy": row[14], "Film-Noir": row[15],
                            "Horror": row[16], "Musical": row[17],
                                "Mystery": row[18], "Romance": row[19], "Sci-Fi": row[20], "
                                    Thriller": row[21], "War": row[22],
                                        "Western": row[23]}
    list_json.append(insert_row)
    file_json = open("Items.json", "w")
35 json.dump(list_json, file_json)
    file_open.close()
    file_json.close()

40 def read_users_data_table():
    path_data_set = "/home/msiddique/WSDL_Work/WebScience/Assignment6/ml-100k/"
    file_open = open(path_data_set + "u.data", "r")
    list_json = []
    reader = csv.reader(file_open, delimiter="\t")
45 for row in reader:
    insert_row = {"user_id": row[0], "item_id": row[1], "rating": row[2], "
        timestamp": row[3]}
    list_json.append(insert_row)
    file_json = open("Data.json", "w")
    json.dump(list_json, file_json)
50 file_open.close()
    file_json.close()

# read_users_data_table()
55 # read_items_table()
read_users_table()
```

```
import json
```

```
def find_closest_user(age, occupation):
5 file_json = open("Users.json", "r")
    list_json = json.load(file_json)
    list_gender = []
```

```
10     for row in list_json:
11         if row["gender"] == "M":
12             list_gender.append(row)
13     list_my_age = []
14     for row in list_gender:
15         if age == int(row["age"]):
16             list_my_age.append(row)
17     list_occupation = []
18     for row in list_my_age:
19         if row["occupation"] == occupation:
20             list_occupation.append(row)
21     if len(list_occupation) < 3:
22         condition_break = False
23         for row in list_my_age:
24             for inner_row in list_occupation:
25                 if row["occupation"] != inner_row["occupation"]:
26                     list_occupation.append(row)
27                     condition_break = True
28                     break
29             if condition_break:
30                 break
31     if len(list_occupation) < 3:
32         condition_break = False
33         for row in list_gender:
34             for inner_row in list_occupation:
35                 if row["user"] != inner_row["user"]:
36                     list_occupation.append(row)
37                     condition_break = True
38                     break
39             if condition_break:
40                 break
41     file_json.close()
42     return list_occupation

43
44 def find_list_of_movie_ratings(age, occupation):
45     list_users = find_closest_user(age, occupation)
46     list_user_id = []
47     for row in list_users:
48         list_user_id.append(row["user"])
49     file_json = open("Data.json", "r")
50     list_json = json.load(file_json)
51     list_favorites = []
52     list_least_favourite = []
53     for user in list_user_id:
54         list_item_id = []
55         list_ratings = []
56         for row in list_json:
57             if row["user_id"] == user:
58                 list_item_id.append(row["item_id"])
59                 list_ratings.append(int(row["rating"]))
60         list_ratings, list_item_id = zip(*sorted(zip(list_ratings, list_item_id)))
61         list_favorites.append(list_item_id[-1])
```

```

        list_least_favourite.append(list_item_id[0])
    file_json.close()
    file_json = open("Items.json", "r")
    list_json = json.load(file_json)
65    list_least_favourite_movies = []
    list_favorites_movies = []
    for item_id in list_favorites:
        for row in list_json:
            if row["movie id"] == item_id:
70                list_favorites_movies.append(row["movie title"])
    for item_id in list_least_favourite:
        for row in list_json:
            if row["movie id"] == item_id:
75                list_least_favourite_movies.append(row["movie title"])

    for i in range(0, len(list_user_id)):
        print("For user id:" + list_user_id[i])
        print("Favorite:" + list_favorites_movies[i])
        print("Favorite movie:" + list_favorites[i])
80        print("Least favorite:" + list_least_favourite[i])
        print("Least favorite movie:" + list_least_favourite_movies[i])

find_list_of_movie_ratings(30, 'student')

```

```

/home/msiddique/WSDL_Work/PolitoWoops_Analysis/venv/bin/python /home/msiddique/WSDL_Work/WebScience/Assignment6/FindClosestUser.py
For user id:774
Favorite:Turbulence (1997)
Favorite movie:986
Least favorite:174
Least favorite movie:Raiders of the Lost Ark (1981)
For user id:869
Favorite:One Fine Day (1996)
Favorite movie:815
Least favorite:118
Least favorite movie:Twister (1996)
For user id:17
Favorite:City of Lost Children, The (1995)
Favorite movie:919
Least favorite:1
Least favorite movie:Toy Story (1995)
Process finished with exit code 0

```

Figure 1: Choice of most similar users to my profile

Problem 2

Which 5 users are most correlated to the substitute you? Which 5 users are least correlated (i.e., negative correlation)?

SOLUTION

For finding correlated users, I created a json in the format shown of critics as shown in the sample recommendation code provided with the assignment. I used the json object to supply it to the function *sim_pearson* from the sample recommendation code with the substituted me and the user ids to generate the correlation values. Figure 2 shows the top 5 least and most correlated users.

```
'''
```

```

/home/msiddique/WSDL_Work/PolitoWoops_Analysis/venv/bin/python /home/msiddique/WSDL_Work/WebScience/Assignment6/CorrelatedUsers.py
Top 5 correlated users
('909', '127', '578', '509', '510')
Last 5 correlated users
('309', '681', '799', '803', '857')
Top 5 recommendation
[(5.0, '814'), (5.0, '1599'), (5.0, '1536'), (5.0, '1467'), (5.0, '1293')]
Last 5 recommendation
[(1.0, '1325'), (1.0, '1321'), (1.0, '1320'), (1.0, '1309'), (1.0, '1308')]
Top 5 correlated movies for my Least favourite movie
('1442', '1558', '1095', '1173', '1249')
Last 5 correlated movies for my least favourite movie
('75', '1465', '113', '119', '1298')
Top 5 correlated movies for my favourite movie
('1508', '883', '1038', '1085', '1102')
Last 5 correlated movies for my favourite movie
('970', '981', '1475', '1015', '1223')
Process finished with exit code 0

```

Figure 2: Top 5 most and least correlated users for substitute me

```

Create a json file of users to mimic the critic dictionary from Recommendation.py
'''
5
def create_critics_json():
    file_json = open("Data.json", "r")
    list_json = json.load(file_json)
    critic={}
10    for row in list_json:
        if row["user_id"] in critic.keys():
            critic[row["user_id"]].update({row["item_id"]: int(row["rating"])})
        else:
            critic[row["user_id"]] = {row["item_id"]: int(row["rating"])}
15    file_json.close()
    file_correlated_users = open("Critic.json", "w")
    json.dump(critic, file_correlated_users)
    file_correlated_users.close()
20
'''
Use Recommendation.py function to calculate correlation and get top 5 and last five
'''
25
def find_correlated_users(user_id):
    create_critics_json()
    file_json = open("Critic.json", "r")
    list_json = json.load(file_json)
30    list_correlation = []
    list_user = []
    for user in list_json:
        list_user.append(user)
        list_correlation.append(sim_pearson(list_json, user_id, user))
35    file_json.close()
    list_correlation, list_user = zip(*sorted(zip(list_correlation, list_user)))
    print("Top 5 correlated users")
    print(list_user[:5])
    print("Last 5 correlated users")
40    print(list_user[-5:])
    get_recommended_list(user_id, list_json)

```



```
choose_real_you(list_json)
```

Problem 3

Compute ratings for all the films that the substitute you have not seen. Provide a list of the top 5 recommendations for films that the substitute you should see. Provide a list of the bottom 5 recommendations (i.e., films the substitute you is almost certain to hate).

SOLUTION

I used the *get_recommendation* function from the recommendation sample code with the json data created from the previous problem of finding correlation users with substitute me user id to get the top 5 most and meast favourite movies. Figure 3 shows the resultsfor recommendation results.

```
/home/msiddique/WSDL_Work/Politoowoops_Analysis/venv/bin/python /home/msiddique/WSDL_Work/WebScience/Assignment6/CorrelatedUsers.py
Top 5 correlated users
('909', '127', '578', '509', '510')
Last 5 correlated users
('309', '681', '799', '803', '857')
Top 5 recommendation
[(5.0, '814'), (5.0, '1599'), (5.0, '1536'), (5.0, '1467'), (5.0, '1293')]
Last 5 recommendation
[(1.0, '1325'), (1.0, '1321'), (1.0, '1320'), (1.0, '1309'), (1.0, '1308')]
Top 5 correlated movies for my Least favourite movie
('1442', '1558', '1095', '1173', '1249')
Last 5 correlated movies for my least favourite movie
('75', '1465', '113', '119', '1298')
Top 5 correlated movies for my favourite movie
('1508', '883', '1038', '1085', '1102')
Last 5 correlated movies for my favourite movie
('970', '981', '1475', '1015', '1223')
Process finished with exit code 0
```

Figure 3: Top 5 favorite and least favorite recommendation for sustitute me user

```
'''
Function to get recommendation list for user
'''

5 def get_recommended_list(user_id, list_json):
    list_recommendation = getRecommendations(list_json, user_id)
    print("Top 5 recommendation")
    print(list_recommendation[:5])
10 print("Last 5 recommendation")
    print(list_recommendation[-5:])
```

Problem 4

Choose your (the real you, not the substitute you) favorite and least favorite film from the data. For each film, generate a list of the top 5 most correlated and bottom 5 least correlated films. Based on your knowledge of the resulting films, do you agree with the results? In other words, do you personally like / dislike the resulting films?

SOLUTION

I added my favourite movie as movie id 161 (Top Gun) and least favorite movie as 162 (On Golden Pond).

I added these entries into the json that I created in the correlation users case and transformed the json to be now listed by movie id rather than user id by using the function *transformPrefs* from the sample recommendation code and ran the *sim_pearson* function to get my list of most and least correlated top 5 movies for my favourite and least favourite movies. Figure 4 shows the movie correlation results.

The results for the correlation values did not make a lot of sense to me, as I had not watched many of the movies listed in the dataset. But looking at their genres and synopsis from IMDB, I realised the correlation results were very similar to my choices.

```
/home/msiddique/WSOL_Work/PolitoWoops_Analysis/venv/bin/python /home/msiddique/WSOL_Work/WebScience/Assignment6/CorrelatedUsers.py
Top 5 correlated users
('909', '127', '578', '509', '510')
Last 5 correlated users
('309', '681', '799', '803', '857')
Top 5 recommendation
[(5.0, '814'), (5.0, '1599'), (5.0, '1536'), (5.0, '1467'), (5.0, '1293')]
Last 5 recommendation
[(1.0, '1325'), (1.0, '1321'), (1.0, '1320'), (1.0, '1309'), (1.0, '1308')]
Top 5 correlated movies for my Least favourite movie
('1442', '1558', '1095', '1173', '1249')
Last 5 correlated movies for my least favourite movie
('75', '1465', '113', '119', '1298')
Top 5 correlated movies for my favourite movie
('1508', '883', '1038', '1085', '1102')
Last 5 correlated movies for my favourite movie
('970', '981', '1475', '1015', '1223')

Process finished with exit code 0
```

Figure 4: Movie Correlation for top 5 most and least correlated movies

```
def choose_real_you(list_json):
    list_json.update({"Me":{"161": 5, "162": 1}})
    list_transformed = transformPrefs(list_json)
    list_correlation_161 = []
    list_correlation_162 = []
    list_movies_161 = []
    list_movies_162 = []
    for movies in list_transformed:
        list_movies_161.append(movies)
        list_movies_162.append(movies)
        list_correlation_161.append(sim_pearson(list_transformed, "161", movies))
        list_correlation_162.append(sim_pearson(list_transformed, "162", movies))
    list_correlation_161, list_movies_161 = zip(*sorted(zip(list_correlation_161,
        list_movies_161)))
    list_correlation_162, list_movies_162 = zip(*sorted(zip(list_correlation_162,
        list_movies_162)))
    print("Top 5 correlated movies for my Least favourite movie")
    print(list_movies_162[:5])
    print("Last 5 correlated movies for my least favourite movie")
    print(list_movies_162[-5:])
    print("Top 5 correlated movies for my favourite movie")
    print(list_movies_161[:5])
    print("Last 5 correlated movies for my favourite movie")
    print(list_movies_161[-5:])
```