

# INTRODUCTION TO WEB SCIENCES: Assignment 7

Manoj Chandra Kompalli

31 March 2016

# Contents

<b>1</b>	<b>Question 1:</b>	<b>2</b>
1.1	Approach . . . . .	2
1.2	Code Listing . . . . .	3
1.2.1	substitute.py to generate favorite and least favorite movies . . . . .	3
1.3	Generated output on shell . . . . .	7
<b>2</b>	<b>Question 2:</b>	<b>8</b>
2.1	Approach . . . . .	8
2.2	Code Listing . . . . .	8
2.2.1	correlation.py . . . . .	8
2.3	Output . . . . .	9
2.3.1	Response showing most and least correlated users to my substitute . . . . .	9
<b>3</b>	<b>Question 3:</b>	<b>10</b>
3.1	Approach . . . . .	10
3.2	Code Listing . . . . .	10
3.2.1	recsub.py . . . . .	10
3.3	Output . . . . .	12
3.3.1	Top movies along with their ratings . . . . .	12
<b>4</b>	<b>Question 4:</b>	<b>12</b>
4.1	Approach . . . . .	12
4.2	Code Listing . . . . .	13
4.2.1	myfav.py . . . . .	13
4.3	Output . . . . .	14
4.3.1	Response showing Pearson's Coefficient along with correlated movies for Shaw-shank Redemption and A Smile like yours . . . . .	14

# 1 Question 1:

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

- what are their top 3 favorite films? - 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".

## 1.1 Approach

Here my task was to find out a substitute of me.

- I have used the u.user file to collect the user data.
- I extracted all the users matching with my age, gender, occupation
- I extracted the top 3 from my list
- For each user, I stored the itemids and rating from u.data file in a dictionary
- I have sorted the dictionary based on the rating
- I have extracted the top three most rated and least rated three items
- I have extracted the movie names for the items using u.items
- I have observed that user 33 has two of my favourite movies like Titanic and Tomorrow never dies.
- Therefore, I chose user 33 to be my substitute

## 1.2 Code Listing

### 1.2.1 substitute.py to generate favorite and least favorite movies

```
1 import collections
2 import sys
3 from itertools import groupby
4 favmovies=[]
5 itemsdict={}
6 itemsdict2={}
7 itemsdict3={}
8 leastdict={}
9 leastdict2={}
10 leastdict3={}
11 def userMovies(movilefile):
12     c=0
13
14
15
16
17     print 'top movies for user 33'
18     print 'itemid rating'
19     for line in movilefile:
20         (user_id, item_id, rating, timestamp) = line.split('\t')
21
22         rating = float(rating)
23         try:
24             if (user_id=='33'):
25                 itemsdict[item_id].append(rating)
26
27         except KeyError:
28             itemsdict[item_id] = list()
29             itemsdict[item_id].append(rating)
30     for key, value in sorted(itemsdict.iteritems(), reverse=True, key=lambda (k,v): (v,k)):
31         c=c+1
32         if (c<4):
33             print "%s: %s" % (key, value)
34             favmovies.append(key)
35
36     c=0
37     movilefile.seek(0)
38     print 'least favourite movies for user 33'
39     print 'itemid rating'
40     for line in movilefile:
41         (user_id, item_id, rating, timestamp) = line.split('\t')
42
43         rating = float(rating)
44         try:
45             if (user_id=='33'):
46                 leastdict[item_id].append(rating)
47
48         except KeyError:
49             leastdict[item_id] = list()
50             leastdict[item_id].append(rating)
51     for key, value in sorted(leastdict.iteritems(), key=lambda (k,v): (v,k)):
52         c=c+1
```

```

53     if(c<4):
54         print "%s: %s" % (key, value)
55
56 c=0
57 moviefile.seek(0)
58 print 'top movies for user 37'
59 print 'itemid rating'
60 for line in movilefile:
61     (user_id, item_id, rating, timestamp) = line.split('\t')
62
63     rating = float(rating)
64     try:
65         if(user_id=='37'):
66             itemsdict2[item_id].append(rating)
67
68     except KeyError:
69         itemsdict2[item_id] = list()
70         itemsdict2[item_id].append(rating)
71 for key, value in sorted(itemsdict2.iteritems(), reverse=True, key=lambda (k,v): (v,k)
72 )):
73     c=c+1
74     if(c<4):
75         print "%s: %s" % (key, value)
76         favmovies.append(key)
77
78 c=0
79 moviefile.seek(0)
80 print 'least favourite movies for user 37'
81 print 'itemid rating'
82 for line in movilefile:
83     (user_id, item_id, rating, timestamp) = line.split('\t')
84
85     rating = float(rating)
86     try:
87         if(user_id=='37'):
88             leastdict2[item_id].append(rating)
89
90     except KeyError:
91         leastdict2[item_id] = list()
92         leastdict2[item_id].append(rating)
93 for key, value in sorted(leastdict2.iteritems(), key=lambda (k,v): (v,k)):
94     c=c+1
95     if(c<4):
96         print "%s: %s" % (key, value)
97
98 c=0
99 moviefile.seek(0)
100 print 'top movies for user 66'
101 print 'itemid rating'
102 for line in movilefile:
103     (user_id, item_id, rating, timestamp) = line.split('\t')
104
105     rating = float(rating)
106     try:
107         if(user_id=='66'):
108             itemsdict3[item_id].append(rating)
109
110     except KeyError:

```

```

109     itemsdict3[item_id] = list()
110     itemsdict3[item_id].append( rating )
111 for key, value in sorted(itemsdict3.iteritems(), reverse=True, key=lambda (k,v): (v,k
112 )):
113     c=c+1
114     if(c<4):
115         print "%s: %s" % (key, value)
116         favmovies.append(key)
117 c=0
118 moviefile.seek(0)
119 print 'least favourite movies for user 66'
120 print 'itemid rating'
121 for line in movilefile:
122     (user_id, item_id, rating, timestamp ) = line.split('\t')
123
124     rating = float(rating)
125     try:
126         if(user_id=='66'):
127             leastdict3[item_id].append(rating)
128
129     except KeyError:
130         leastdict3[item_id] = list()
131         leastdict3[item_id].append( rating )
132 for key, value in sorted(leastdict3.iteritems(), key=lambda (k,v): (v,k)):
133     c=c+1
134     if(c<4):
135         print "%s: %s" % (key, value)
136
137 c=0
138 moviefile.seek(0)
139 def getUsers( userfile ):
140     count=0
141     for line in userfile:
142         line=line.split("|")
143         id=line[0]
144         age=line[1]
145         gender=line[2]
146         occupation=line[3]
147
148         if(age=='23' and gender=='M' and occupation=='student' and count<3):
149             count=count+1
150             print 'top 3 substitutes'
151             print 'id'+ ' ' + 'count'+ ' ' + 'occupation'+ ' ' + 'gender'+ ' ' + 'age'
152             print id, count, occupation, gender, age
153
154 def getTitles(imdbfile):
155     ct=0
156     namesdict = {}
157     print 'top movies of my substitutes'
158     for line in imdbfile:
159         # print line
160         line = line.split("|")
161         id = line[0]
162         title = line[1]
163
164         namesdict[ id ] = title

```

```

165     for key,value in namesdict.items():
166         try:
167             if key==favmovies[ct] and ct<9:
168                 ct=ct+1
169                 print(key,value)
170         except:
171             pass
172     print 'my final substitute is user 33 because i like 2 of his movies titanic and
173         tomorrow never dies'
174
175
176
177
178
179 if __name__=='__main__':
180     userfile=open("u.user","r")
181     moviefile=open("u.data","r")
182     imdbfile=open("u.item","r")
183     users=getUsers(userfile)
184     movies=userMovies(moviefile)
185     titles=getTitles(imdbfile)

```

## 1.3 Generated output on shell

```
sirius:~/wsprograms/a7new/a7/q1> python substitute.py
top movies for user 33
itemid rating
313: [5.0]
751: [4.0]
682: [4.0]
least favourite movies for user 33
itemid rating
245: [3.0]
294: [3.0]
307: [3.0]
top movies for user 37
itemid rating
68: [5.0]
62: [5.0]
597: [5.0]
least favourite movies for user 37
itemid rating
82: [1.0]
118: [2.0]
121: [2.0]
top movies for user 66
itemid rating
742: [5.0]
50: [5.0]
471: [5.0]
least favourite movies for user 66
itemid rating
21: [1.0]
286: [1.0]
877: [1.0]
top movies of my substitutes
('313', 'Titanic (1997)')
('751', 'Tomorrow Never Dies (1997)')
('682', 'I Know What You Did Last Summer (1997)')
('68', 'Crow, The (1994)')
('62', 'Stargate (1994)')
('597', 'Eraser (1996)')
('742', 'Ransom (1996)')
('50', 'Star Wars (1977)')
('471', 'Courage Under Fire (1996)')
my final substitute is user 33 because i like 2 of his movies titanic and tomorrow never dies
sirius:~/wsprograms/a7new/a7/q1> █
```

Figure 1: Top and least rated movies and their ratings of user 33



## 2 Question 2:

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

### 2.1 Approach

Here, I used the recommendations.py file mentioned in the book Programming Collective Intelligence

- I have used topMatches() method to extract the users correlated with my substitute user(user 33).
- I have stored the results in a dictionary and extracted the top five users from it for the users with most correlation
- I have extracted the bottom five users from the same dictionary to get the bottom five for the users with least correlation

### 2.2 Code Listing

#### 2.2.1 correlation.py

```
1 #!/usr/local/bin/python
2
3 import sys
4 import recommendations
5 import operator
6 def getRatings(ratingsfile):
7     itemsdict = {}
8     count=0
9     for line in ratingsfile:
10         (user_id, item_id, rating, timestamp) = line.split('\t')
11         if user_id in itemsdict:
12             itemsdict[user_id][item_id] = float(rating)
13
14         else:
15             count=count+1
16             itemsdict[user_id]={}
17 output= recommendations.topMatches(itemsdict, '33',n=count)
18 top = output[:5]
19 print 'top 5 correlated users to user 33 are'
20 print top
21 bot = output[-5:]
22 print 'bottom 5 correlated users to user 33 are'
23 print bot
24 if __name__ == '__main__':
25
26
27     ratingsfile = open("../ml-100k/u.data", "r")
28     ratings      = getRatings(ratingsfile)
```

## 2.3 Output

### 2.3.1 Response showing most and least correlated users to my substitute

```
sirius:~/wsprograms/a7new/a7/q2> python correlation.py
top 5 correlated users to user 33 are
[(1.0000000000000004, '890'), (1.00000000000000018, '937'), (1.00000000000000018, '718'), (1.0, '95'), (1.0, '935')]
bottom 5 correlated users to user 33 are
[(-1.0, '136'), (-1.0, '103'), (-1.0000000000000001, '736'), (-1.00000000000000018, '657'), (-1.00000000000000027, '797'
)]
sirius:~/wsprograms/a7new/a7/q2> █
```

Figure 2: Pearson Coefficients and Correlated users to user 33

## 3 Question 3:

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).

### 3.1 Approach

My substitute you has seen and liked a fair bit of movies. Based on this dataset I have to find the most recommended movies and least recommended movies for my substitute.

- I used the `getRecommendations` method from `recommendations.py` which returned me all the recommended movies for a particular user.
- I fetched the movie titles and their respective ratings for all the resulting movie titles
- The resulting movies are sorted based on their ratings

### 3.2 Code Listing

#### 3.2.1 recsub.py

```
1 import recommendations
2 import sys
3 from itertools import groupby
4 itemsdict = {}
5 def getTitles(imdbfile, movie_list):
6     ct=0
7     namesdict = {}
8     movie_id = []
9     ratings=[]
10    for i in movie_list:
11        movie_id.append(i[1])
12        ratings.append(i[0])
13    movie_titles = []
14
15    for line in imdbfile:
16        # print line
17        line = line.split("|")
18        id = line[0]
19        title = line[1]
20
21        namesdict[ id ] = title
22        if id in movie_id:
23            movie_titles.append(title)
24
25    return movie_titles, ratings
26
27 if __name__ == '__main__':
28     moviefile = open("u.data", "r")
29
30     pref_dic = {}
31     rec=[]
32     for line in moviefile:
```

```

33     (user_id, item_id, rating, timestamp) = line.split('\t')
34     if user_id in pref_dic:
35         pref_dic[user_id][item_id] = rating
36     else:
37         pref_dic[user_id] = {}
38     result = recommendations.getRecommendations(pref_dic, '33')
39     rec.append(result)
40
41
42
43
44
45
46     bottom5= result[-5:]
47
48     top5 = result[:5]
49     imdbfile=open("u.item", "r")
50
51     top_titles, rate=getTitles(imdbfile, top5)
52     print "Top Movies are"
53     for i in range(0, len(top_titles)):
54         print top_titles[i], rate[i]
55
56     imdbfile=open("u.item", "r")
57     bottom_titles, rate = getTitles(imdbfile, bottom5)
58     print "Bottom Movies are"
59     for j in range(0, len(bottom_titles)):
60         print bottom_titles[j], rate[j]

```

## 3.3 Output

### 3.3.1 Top movies along with their ratings

```
sirius:~/wsprograms/a7new/a7/q3> python recsub.py
Top Movies are
Great Day in Harlem, A (1994) 5.0
Two or Three Things I Know About Her (1966) 5.0
Hearts and Minds (1996) 5.0
Prefontaine (1997) 5.0
Santa with Muscles (1996) 5.0
Bottom Movies are
Theodore Rex (1995) 1.0
Country Life (1994) 1.0
New York Cop (1996) 1.0
Catwalk (1995) 1.0
Homage (1995) 1.0
sirius:~/wsprograms/a7new/a7/q3> █
```

Figure 3: Recommended movies sorted based on ratings

## 4 Question 4:

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?

### 4.1 Approach

My favorite movie is Shawshank Redemption by a long way. Although there are many movies which belong to least favorite movie, I chose Smile like yours. Its rated 4.2 in IMDB right now. Now, my next job is to get the most correlated movies and least correlated movies for both my favorite and least favorite movies.

- To show correlation, again I have used Pearsons Coefficients.
- Some of the key functions I used from recommendations.py to extract the data are loadMovieLens which loads the movie lens data, topMatches which sorts the dictionary containing movie items and pearson coefficient.
- TopMatches again calls a method “simpearson” which finds the preferences and ultimately Pearson score.
- A Pearson score closer to 1 shows highly correlated movies and score close to -1 shows least correlated movies.
- I have calculated all twenty scores for my favorite and least favorite films.
- Unfortunately though I cannot comment on whether I like or dislike the resulting movies because I have hardly watched any of them

## 4.2 Code Listing

### 4.2.1 myfav.py

```
1 #!/usr/local/bin/python
2 import sys
3 import recommendationsedited
4
5 if __name__ == '__main__':
6
7     outlier = sys.argv[1]
8     movie = 'Shawshank Redemption, The (1994)'
9     movie2='Smile Like Yours, A (1997)'
10
11     count = 5
12
13     if outlier not in [ "top5", "bottom5" ]:
14         print "argument 1 must be either top5 or bottom5"
15         sys.exit(1)
16
17
18     prefs = recommendationsedited.loadMovieLens( '../data' )
19     itemPrefs = recommendationsedited.transformPrefs( prefs )
20     results = recommendationsedited.topMatches( itemPrefs, movie, 2000 )
21     results2= recommendationsedited.topMatches( itemPrefs, movie2, 2000 )
22
23     if outlier == "top5":
24         print 'for favourite movie Shawshank Redemption most correlated movies are'
25         for i in results[0:count]:
26
27             print i[0], i[1]
28         print 'for least favourite movie A Smile like yours most correlated movies are'
29         for i in results2[0:count]:
30
31             print i[0], i[1]
32
33
34     if outlier == "bottom5":
35
36
37         results.reverse()
38         results2.reverse()
39         print 'for favourite movie Shawshank Redemption least correlated movies are'
40         for i in results[0:count]:
41             #print i[0], i[1]
42         print i
43         print 'for least favourite movie A Smile like yours bottom5 least correlated'
44         print 'movies are'
45         for i in results2[0:count]:
46             #print i[0], i[1]
47         print i
```

## 4.3 Output

### 4.3.1 Response showing Pearson's Coefficient along with correlated movies for Shawshank Redemption and A Smile like yours

```
sirius:~/wsprograms/a7new/a7/q4> python myfav.py "myfav.py"
argument 1 must be either top5 or bottom5
sirius:~/wsprograms/a7new/a7/q4> python myfav.py "top5"
for favourite movie Shawshank Redemption most correlated movies are
1.0 Penny Serenade (1941)
1.0 Newton Boys, The (1998)
1.0 Oscar & Lucinda (1997)
1.0 Wedding Gift, The (1994)
1.0 Search for One-eye Jimmy, The (1996)
for least favourite movie A Smile like yours most correlated movies are
1.0 City of Angels (1998)
1.0 Truth About Cats & Dogs, The (1996)
1.0 Silence of the Lambs, The (1991)
1.0 Othello (1995)
1.0 Fargo (1996)
sirius:~/wsprograms/a7new/a7/q4> python myfav.py "bottom5"
for favourite movie Shawshank Redemption least correlated movies are
(-1.0000000000000004, 'Clean Slate (Coup de Torchon) (1981)')
(-1.0, '1-900 (1994)')
(-1.0, 'American Dream (1990)')
(-1.0, 'Collectionneuse, La (1967)')
(-1.0, 'Colonel Chabert, Le (1994)')
for least favourite movie A Smile like yours bottom5 least correlated movies are
(-1.0000000000000007, 'Andre (1994)')
(-1.0, 'Flubber (1997)')
(-1.0, 'Nightwatch (1997)')
(-1.0, 'Vegas Vacation (1997)')
(-1.0, 'While You Were Sleeping (1995)')
sirius:~/wsprograms/a7new/a7/q4> █
```

Figure 4: Pearson's Coefficient and movie name

## References

- [1] Movie Lens dataset <http://grouplens.org/datasets/movielens/100k/>,
- [2] Programming Collective Intelligence <http://forum.myquant.cn>,
- [3] Recommendations.py <https://github.com/arthur-e/Programming-Collective-Intelligence/blob/master/chapter2/recommendations.py>

□