# HW 7 - Recommendation Systems
Adeniran Adeniyi
DUE Sunday, April 11, 2021 before 11:59pm

# Q1

*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?

- what are their bottom 3 (least favorite) films?

## Answer

```python
# -*- coding: utf-8 -*-
"""
Created on Wed Apr  7 20:09:11 2021

@author: aadeniran
"""
import pandas as pd

def closetToMeAgeGendOccu(item,data):

    #load movies data into variables
    movies = {}
    for line in open(str(item)):
        (id, title) = line.split('|')[0:2]
        movies[id] = title
    # Load data
    prefs = {}
    for line in open(str(data)):
        (user, movieid, rating, ts) = line.split('\t')
        prefs.setdefault(user, {})
        prefs[user][movies[movieid]] = float(rating)
    #Load in users
    users = []
    for line in open('movie/u.user'):
        (user, age, gender, occupation, zipcode) = line.split('|')
        dictA = {'user': user, 'age':age, 'gender': gender, 'occupation' :occupation }
        users.append(dictA)
```

```
28
29      count = 1
30      #hold 3 similar user in term of age gender and category
31      similarUser =[]
32      #Find all users that have things common with me
33      for a in users:
34          """
35          Here i parsed in my age gender and occupation
36          to find users of my age gender and occupation category.
37          I take only 3 users
38          """
39          if a['age'] == '24' and a['gender'] == 'M' and a['occupation']
    == 'student' and count < 4 :
40              similarUser.append(a['user'])
41              count += 1
42
43       #print user movie likes
44       for a in similarUser:
45          """
46          Based on similar users like me,
47          I obtained these similar users top 3 favorite films
48          """
49          if a in prefs:
50              #this gets the movie names, and the movie rating per user
51              p = pd.DataFrame(list(prefs[a].items()),columns=["Movie
    Names","Rating"])
52              #this sort the users movie list based on rating in
    decending order
53              p.sort_values("Rating",ascending=False, inplace=True)
54              #gets the top3 films
55              top3 = p.head(3)
56              #gets the bottom 3 films
57              bottom3 = p.tail(3)
58              #print("Top 3 films for User {}:".format(a))
59              #print(top3)
60              #print("\n\n")
61              #print("Bottom 3 films for User {}:".format(a))
62              #print(bottom3)
63              #print("\n")
64              top3.to_csv("Q1/"+a+"top3.csv",index=False)
65              bottom3.to_csv("Q1/"+a+"bottom.csv",index=False)
66
67
68
69 """
70 Question 1 solution
71 """
```

```
72 closetToMeAgeGendOccu('movie/u.item','movie/u.data')
```

**Listing 1:** movieLen.py

**Table 1:** Top 3 movies for user 73

| Movie Names | Ratings |
|---|---|
| Three Colors: Red (1994) | 5.0 |
| Godfather: Part II, The (1974) | 5.0 |
| 2001: A Space Odyssey (1968) | 5.0 |

**Table 2:** Bottom 3 movies for user 73

| Movie Names | Ratings |
|---|---|
| Saint, The (1997) | 2.0 |
| Home Alone 3 (1997) | 1.0 |
| Home Alone (1990) | 1.0 |

**Table 3:** Top 3 movies for user 301

| Movie Names | Ratings |
|---|---|
| It's a Wonderful Life (1946) | 5.0 |
| Empire Strikes Back, The (1980) | 5.0 |
| Star Wars (1977) | 5.0 |

**Table 4:** Bottom 3 movies for user 301

| Movie Names | Ratings |
|---|---|
| Ready to Wear (Pret-A-Porter) (1994) | 1.0 |
| Dirty Dancing (1987) | 1.0 |
| Natural Born Killers (1994) | 1.0 |

**Table 5:** Top 3 movies for user 369

| Movie Names | Ratings |
|---|---|
| Dead Poets Society (1989) | 5.0 |
| Return of the Jedi (1983) | 5.0 |
| Wallace & Gromit: The Best of Aardman Animation (1996) | 5.0 |

**Table 6:** Bottom 3 movies for user 369

| Movie Names | Ratings |
|---|---|
| Beautician and the Beast, The (1997) | 3.0 |
| How to Be a Player (1997) | 2.0 |
| Booty Call (1997) | 2.0 |

## Discussion

*I created a function called closetToMeAgeGendOccu(),in line 9. This function takes two parameters of the location of u.item and u.data path. u.user was not going through as a parameter so i parsed it directly. From lines 11-27 I read in the datas from each file and store it into dictionary varaibles. In Line 29 -41, I found users that were just as similar to me in age gender and occupation. I limited this search to find only just 3 similar users to me on line 35. These users where stored in a list variable called similarUser. line 31 is the variable declaration, while line 40 appends new users to the list in the created variable*

*Lines 44 to 65, process the outcomes for top 3 movies and bottom 3 movies for each users that I found. Using a for loop to retrieve each user, parse this result to ensure that the said user is in prefs dictionary. If present, retrieve the row of the dictionary of that particular user. Using a pandas dataFrame in Line 51, parse a list of the row item, so that we can easily sort the data based on values in Rating column in descending order. Retrieve the first three data using pandas function .head(3) and the last three data using .tail(3). Save each result as a pandas dataFrame and convert the result to a saved csv file in Q1.*

*My best substitute user is 301, the movies he hates I hate and the movies he enjoys I find them interesting too*

# Q2

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

## Answer

```python
1  # -*- coding: utf-8 -*-
2  """
3  Created on Thu Apr  8 00:46:12 2021
4
5  @author: aadeniran
6  """
7  from math import sqrt
8
9  def sim_pearson(prefs, p1, p2):
10     """
11     Returns the Pearson correlation coefficient for p1 and p2.
12     """
13
14 # Get the list of mutually rated items
15     si = {}
16     for item in prefs[p1]:
17         if item in prefs[p2]:
18             si[item] = 1
19
20 # If they are no ratings in common, return 0
21     if len(si) == 0:
22         return 0
23
24 # Sum calculations
25     n = len(si)
26
27  # Sums of all the preferences
28     sum1 = sum([prefs[p1][it] for it in si])
29     sum2 = sum([prefs[p2][it] for it in si])
30
31 # Sums of the squares
32     sum1Sq = sum([pow(prefs[p1][it], 2) for it in si])
33     sum2Sq = sum([pow(prefs[p2][it], 2) for it in si])
34
35  # Sum of the products
36     pSum = sum([prefs[p1][it] * prefs[p2][it] for it in si])
37
38  # Calculate r (Pearson score)
```

```
39      num = pSum - sum1 * sum2 / n
40      den = sqrt((sum1Sq - pow(sum1, 2) / n) * (sum2Sq - pow(sum2, 2) / n
41      ))
42      if den == 0:
43          return 0
44      r = num / den
45      return r
46
47  def getRecommendations(prefs, person, similarity=sim_pearson):
48      """
49       Gets recommendations for a person by using a weighted average
50       of every other  u s e r s  rankings
51      """
52      totals = {}
53      simSums = {}
54      for other in prefs:
55          # D o n t  compare me to myself
56          if other == person:
57              continue
58          sim = similarity(prefs, person, other)
59          # Ignore scores of zero or lower
60          if sim <= 0:
61              continue
62          for item in prefs[other]:
63              # Only score movies I  h a v e n t  seen yet
64              if item not in prefs[person] or prefs[person][item] == 0:
65                  # Similarity * Score
66                  totals.setdefault(item, 0)
67                  # The final score is calculated by multiplying each
     item by the
68                  # similarity and adding these products together
69                  totals[item] += prefs[other][item] * sim
70                  # Sum of similarities
71                  simSums.setdefault(item, 0)
72                  simSums[item] += sim
73                  # Create the normalized list
74          rankings = [(total / simSums[item], item) for (item, total) in
75                          totals.items()]
76          # Return the sorted list
77          rankings.sort()
78          rankings.reverse()
79          return rankings
80
81  def topMatches(prefs,person,n=5,similarity=sim_pearson,):
82      '''
83      Returns the best matches for person from the prefs dictionary.
```

```python
84      Number of results and similarity function are optional params.
        bottomMatches
85      '''
86
87      scores = [(similarity(prefs, person, other), other) for other in
        prefs
88              if other != person]
89      scores.sort()
90      scores.reverse()
91      return scores[0:n]
92
93  def bottomMatches(prefs,person,n=5,similarity=sim_pearson,):
94      """
95      Returns the lowest matches for person from the prefs dictionary.
96      Number of results and similarity function are optional params.
97      """
98      scores = [(similarity(prefs, person, other), other) for other in
        prefs
99              if other != person]
100     scores.sort()
101     scores.reverse()
102     return scores[len(scores)-n: len(scores)]
103
104 def transformPrefs(prefs):
105     '''
106     Transform the recommendations into a mapping where persons are
        described
107     with interest scores for a given title e.g. {title: person} instead
        of
108     {person: title}.
109     '''
110
111     result = {}
112     for person in prefs:
113         for item in prefs[person]:
114             result.setdefault(item, {})
115             # Flip item and person
116             result[item][person] = prefs[person][item]
117     return result
118
119 def loadMovieLens():
120     # Get movie titles
121     movies = {}
122     for line in open("movie/u.item"):
123         (id, title) = line.split("|")[0:2]
124         movies[id] = title
125     # Load data
```

```
126      #prefs = {}
127      for line in open("movie/u.data"):
128          (user, movieid, rating, ts) = line.split("\t")
129          prefs.setdefault(user, {})
130          prefs[user][movies[movieid]] = float(rating)
131      #Load in users
132      #users = []
133      for line in open("movie/u.user"):
134          (user, age, gender, occupation, zipcode) = line.split("|")
135          dictA = {"user": user, "age":age, "gender": gender, "occupation
     " :occupation }
136      users.append(dictA)
137
138
139 prefs ={}
140 users =[]
141 loadMovieLens()
142
143 top5 = topMatches(prefs, "301", n = 5, similarity = sim_pearson)
144 bottom5 = bottomMatches(prefs, "928", n = 5, similarity = sim_pearson)
145 """
146 Question 2
147 (1.0000000000000029, '801')
148 (1.0, '845')
149 (1.0, '140')
150 (1.0, '111')
151 (0.999999999999992, '565')
152
153
154 (-1.000000000000004, '760')
155 (-1.000000000000004, '547')
156 (-1.000000000000004, '432')
157 (-1.000000000000004, '317')
158 (-1.000000000000004, '112')
159 """
160
161 print("5 users that  most correlate with my substitute me ")
162 print(*top5, sep='\n')
163 print("\n\n")
164 print("5 users that least correlate with my substitute me ")
165 print(*bottom5, sep='\n')
166
167 """
168 Question 3
169 (5.0, 'Unforgettable (1996)')
170 (5.0, 'Net, The (1995)')
171 (5.0, 'Murder in the First (1995)')
```

```
172  (5.0, 'Murder at 1600 (1997)')
173  (5.0, 'Just Cause (1995)')
174
175
176  (2.0, 'Diabolique (1996)')
177  (2.0, 'Devil in a Blue Dress (1995)')
178  (2.0, 'Breakdown (1997)')
179  (1.0, 'Thinner (1996)')
180  (1.0, 'Albino Alligator (1996)')
181  """
182  recommend = getRecommendations(prefs,"301")
183  print("\nTop 5 movies recommendations for substitute me: ")
184  print(*recommend[0:5], sep="\n")
185  print("\nBottom 5 movies recommendations for substitute me:")
186  print(*recommend[len(recommend)-5: len(recommend)], sep="\n")
187
188
189  """
190  Question 4
191
192  [(1.000000000000004, 'Young Guns II (1990)'), (1.000000000000004, '
          Reality Bites (1994)'), (1.000000000000001, 'Ran (1985)'),
          (1.000000000000001, 'Butch Cassidy and the Sundance Kid (1969)'),
          (1.000000000000009, 'Chinatown (1974)')]
193
194
195
196  [(-1.0, 'Big Night (1996)'), (-1.0, 'Barbarella (1968)'), (-1.0, '
          Another Stakeout (1993)'), (-1.000000000000007, "Fathers' Day
          (1997)"), (-1.000000000000004, "Devil's Own, The (1997)")]
197
198  """
199  print("\n\n")
200  movies = transformPrefs(prefs)
201  mybestmovie = "Vampire in Brooklyn (1995)"
202  topFive = topMatches(movies, mybestmovie)
203
204  print("My best 5 recommended movies")
205  print(topFive)
206
207  print("My worst 5 recommended movies")
208  worstFive = bottomMatches(movies,mybestmovie)
209  print("\n")
210  print(worstFive)
```

**Listing 2:** question2.py

```
5 users that  most correlate with my substitute me
(1.0000000000000029, '801')
(1.0, '845')
(1.0, '140')
(1.0, '111')
(0.99999999999992, '565')




5 users that least correlate with my substitute me
(-1.0000000000000004, '760')
(-1.0000000000000004, '547')
(-1.0000000000000004, '432')
(-1.0000000000000004, '317')
(-1.0000000000000004, '112')
```

**Figure 1:** The resulting out put for Question 2

## Discussion

The Final output for the top 5 users that were most correlated to my substitute me were 801,845,140,111, and 565 with their respective correlating score of 1.0000000000000029,1.0,1.0,1.0 and 0.99999999999992.

For the bottom 5 users that were most correlated to my substitute me were, 760,547,432,317,and 112. Their correlative score was the same, -1.000000000000004.

For the driver lines 139 to 165 generated the resulting output. First I loaded the u.items, u.data and u.user in this function and stored them in a dictionary, dictionary and list variables respectively. I called the topMatches function which I parsed in my substitute me user id in.

## Q3

Compute ratings for all the films that the substitute you has 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).*

## Answer

Reference question2.py in Listing 2

```
Top 5 movies recommendations for substitute me:
(5.0, 'Unforgettable (1996)')
(5.0, 'Net, The (1995)')
(5.0, 'Murder in the First (1995)')
(5.0, 'Murder at 1600 (1997)')
(5.0, 'Just Cause (1995)')

Bottom 5 movies recommendations for substitute me:
(2.0, 'Diabolique (1996)')
(2.0, 'Devil in a Blue Dress (1995)')
(2.0, 'Breakdown (1997)')
(1.0, 'Thinner (1996)')
(1.0, 'Albino Alligator (1996)')
```

**Figure 2:** The resulting output for Question 3

## Discussion

*In line 182 to 186, function getRecommendations(prefs,"301") was to get my substiture user recommended list, I Used \*recommend[0:5] to get first 5 recommendation. This code was gotten from* `https://github.com/arthur-e/Programming-Collective-Intelligence/` `blob/master/chapter2/recommendations.py`
*Top 5 movies recommendations for substitute me, Unforgettable (1996), Net The (1995), Murder in the First (1995), Murder at 1600 (1997), Just Cause (1995)*
*Bottom 5 movies recommendations for substitute me,in line 186 \*recommend[len(recommend)-5: len(recommend)] used it to extract the bottom 5 movie recommendations for substitute me Diabolique (1996), Devil in a Blue Dress (1995), Breakdown (1997), Thinner (1996), Albino Alligator (1996)*

# Q4

*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?*

## Answer

Reference question2.py in Listing 2

```
My best 5 recommended movies
[(1.000000000000004, 'Young Guns II (1990)'),
(1.000000000000004, 'Reality Bites (1994)'),
(1.000000000000001, 'Ran (1985)'),
(1.000000000000001, 'Butch Cassidy and the Sundance
Kid (1969)'), (1.0000000000000009, 'Chinatown
(1974)')]
```

```
My worst 5 recommended movies
[(-1.0, 'Big Night (1996)'), (-1.0, 'Barbarella
(1968)'), (-1.0, 'Another Stakeout (1993)'),
(-1.0000000000000007, "Fathers' Day (1997)"),
(-1.000000000000004, "Devil's Own, The (1997)")]
```

**Figure 3:** The resulting output for Question 4

## Discussion

*I first transformed the prefs dictionary to movies. I then selected my best (Vampire in Brooklyn (1995)) movie from the whole list in u.data. The function topMatches and bottomMatch handled the result of the questions 3*

*Top 5 movies recommend from (Vampire in Brooklyn (1995)  Young Guns II (1990)* `https://www.youtube.com/watch?v=r-FmfxLy7fo` *I like the movie, a good old movie to me.*

*Reality Bites (1994)* `https://www.youtube.com/watch?v=xDYGo0UgIVM` *It seem interesting and I would watch it. Just an average movie for me.*
*Ran (1985)* `https://www.youtube.com/watch?v=YwP_kXyd-Rw` *I love this one also, war action my favorite.*
*Butch Cassidy and the Sundance Kid (1969)* `https://www.youtube.com/watch?v=YdJW2UxvSFQ` *I do love cow boys like movies*
*Chinatown (1974)* `https://www.youtube.com/watch?v=20FfiP7g4tU` *Quite boring but I would enjoy the detective part to movie*

*Worst 5 movies recommended from (Vampire in Brooklyn (1995)*

*Big Night (1996)* `https://www.youtube.com/watch?v=Yd8gK6EgpLM` *looks boring to me because its a chef movie*
*Barbarella (1968)* `https://www.youtube.com/watch?v=M-fJg08wBKw` *A complete no to me because its too old and i dont find it interesting at all*
*Another Stakeout (1993)* `https://www.youtube.com/watch?v=Gpm4lGyOVYc` *I like this one, action, detective kind of movies*
*Fathers' Day(1997)* `https://www.youtube.com/watch?v=xsQfKt08Xlk` *This kind of movies does not capture my interest at all. Too boring.*
*Devil's Own, The (1997)* `Devil\OT1\textquoterightsOwn,The(1997)` *What I cant believe this is a least favourite movie. I completely love the action in this movie just the first second of seeing it. As usual war movies gets me all the time.]*

# References

- `https://github.com/arthur-e/Programming-Collective-Intelligence/blob/master/chapter2/recommendations.py`

- `https://www.example.com/reallyreallyreally-extra-long-URI/`