

INTRODUCTION TO WEB SCIENCES: Assignment 8

Babitha Bokka

15 November 2014

Contents

1	Question 1:	4
1.1	Approach	4
1.2	Source Code	5
1.2.1	averageMovies.py	5
1.3	Input Files	6
1.3.1	u.data	6
1.3.2	u.item	6
1.4	Output Files	7
1.4.1	averageMovies.png	7
1.4.2	averageMoviesMore.txt	8
2	Question 2:	9
2.1	Approach	9
2.2	Source Code	10
2.2.1	mostRatingMovies.py	10
2.3	Input Files	11
2.3.1	u.data	11
2.3.2	u.item	11
2.4	Output Files	12
2.4.1	mostRatingMovies.png	12
3	Question 3:	13
3.1	Approach	13
3.2	Source Code	14
3.2.1	averageMovies.py	14
3.3	Input Files	16
3.3.1	u.data	16
3.3.2	u.item	16
3.3.3	u.user	16
3.4	Output Files	17
3.4.1	averageGenderF.png	17
3.4.2	averageGenderFull.txt	18
4	Question 4:	19
4.1	Approach	19
4.2	Source Code	20
4.2.1	averageMovies.py	20
4.3	Input Files	22
4.3.1	u.data	22
4.3.2	u.item	22
4.3.3	u.user	22
4.4	Output Files	23
4.4.1	averageGenderM.png	23
4.4.2	averageGenderFull.txt	24

5	Question 5:	25
5.1	Approach	25
5.2	Source Code	26
5.2.1	recommendations.py	26
5.2.2	ratingsTopGun.py	27
5.3	Input Files	28
5.3.1	u.data	28
5.3.2	u.item	28
5.4	Output Files	29
5.4.1	ratingsLikeTopGun.txt	29
5.4.2	ratingsLeastTopGun.txt	31
6	Question 6:	33
6.1	Approach	33
6.2	Source Code	34
6.2.1	ratarsMostRatedMovies.py	34
6.3	Input Files	35
6.3.1	u.data	35
6.4	Output Files	35
6.4.1	ratarsMostRatedMovies.png	35
7	Question 7:	36
7.1	Approach	36
7.2	Description of	36
7.3	Source Code	36
7.4	Output Files	36
8	Question 8:	37
8.1	Approach	37
8.2	Description of	37
8.3	Source Code	37
8.4	Output Files	37
9	Question 9:	38
9.1	Approach	38
9.2	Source Code	39
9.2.1	averageOverUnderGender.py	39
9.3	Input Files	41
9.3.1	u.data	41
9.3.2	u.item	41
9.3.3	u.user	41
9.4	Output Files	42
9.4.1	averageOverM.png	42
9.4.2	averageUnderM.png	43
9.4.3	averageOverM.txt	44
9.4.4	averageUnderM.txt	45

10 Question 10:	46
10.1 Approach	46
10.2 Source Code	47
10.2.1 averageOverUnderGender.py	47
10.3 Input Files	49
10.3.1 u.data	49
10.3.2 u.item	49
10.3.3 u.user	49
10.4 Output Files	50
10.4.1 averageOverF.png	50
10.4.2 averageUnderF.png	51
10.4.3 averageOverF.txt	52
10.4.4 averageUnderF.txt	53

1 Question 1:

What 5 movies have the highest average ratings? Show the movies and their ratings sorted by their average ratings.

1.1 Approach

1. To solve this question the straight forward approach is to read u.data file get movie id, corresponding movie ratings and read u.item file to get movie id and corresponding movie name.
2. Extract the movie id and ratings. For each movie id build a dictionary that has all the ratings.
3. Now, for each movie id calculate the average using `sum()` and `len()` functions.
4. For each movie id get the movie name and append it to the dictionary which has average.
5. Now, sort the dictionary and print the top 5 movies with highest average ratings.
6. Program `averageMovies.py` produces the top 5 movies with highest average ratings. The output is based on a condition which can be changed accordingly.
7. Figure 1 shows the results of the program.

1.2 Source Code

1.2.1 averageMovies.py

```
1 #!/usr/bin/env python
2
3 import sys
4
5 def main():
6     # Create a dictionary
7     movie_id_ratings = {}
8     movie_id_avg = {}
9     count = 1
10    # Read the input files.
11    readData = open('/home/bbokka/cs594/A8/u.data', 'r')
12    readItem = open('/home/bbokka/cs594/A8/u.item', 'r')
13    # reading the u.data file for item and ratings.
14    for line in readData:
15        split_input_line = line.strip().split('\t')
16        movie_id = split_input_line[1]
17        movie_rating = float(split_input_line[2])
18        try:
19            movie_id_ratings[movie_id].append(movie_rating)
20        except KeyError:
21            movie_id_ratings[movie_id] = list()
22            movie_id_ratings[movie_id].append(movie_rating)
23    readData.close()
24    # Calculating the average.
25    for key in movie_id_ratings:
26        avg = sum(movie_id_ratings[key])/len(movie_id_ratings[key])
27        movie_id_avg[key] = [float(avg)]
28
29    # Reading the u.item file for movie name.
30    for each_line in readItem:
31        split_each_line = each_line.strip().split('|')
32        movie_item_id = split_each_line[0]
33        movie_name_split = split_each_line[1].split('(1')
34        movie_name = movie_name_split[0]
35
36        movie_id_avg[movie_item_id].append(movie_name)
37    readItem.close()
38    print "*" * 60
39    print "Top 5 movies having the highest average ratings"
40    print "*" * 60
41    print "Movie Name\t\t\t\t\t","Avg Rating"
42    print "-" * 60
43    # sorting the movies from highest to lowest based on the average value.
44    for key, value in sorted(movie_id_avg.items(), key=lambda e: e[1], reverse=True):
45        if (count <= 30):
46            print '{:<50}{:<0.1f} '.format(value[1], value[0])
47            count += 1
48
49 if __name__ == "__main__":
50     try:
51         main()
52     except KeyboardInterrupt:
53         sys.exit(1)
```

1.3 Input Files

1.3.1 u.data

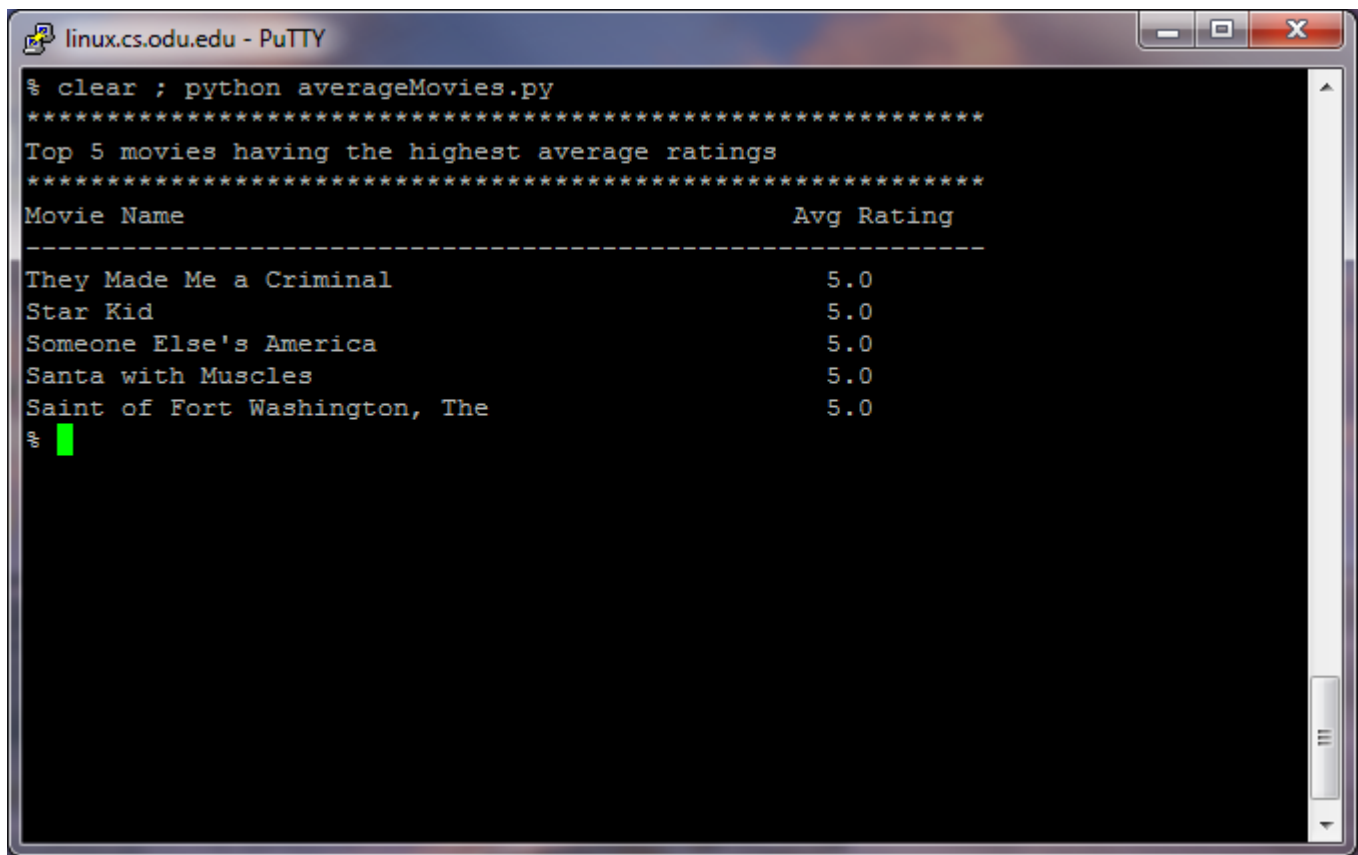
	userid	itemid	rating	timestamp
1				
2				
3	196	242	3	881250949
4	186	302	3	891717742
5	22	377	1	878887116
6	244	51	2	880606923
7	166	346	1	886397596
8	298	474	4	884182806
9	115	265	2	881171488

1.3.2 u.item

movie id	movie title	release date	video release date	IMDb URL	unknown
	Action Adventure Animation Children's Comedy Crime Documentary Drama Fantasy Film-Noir Horror Musical Mystery Romance Sci-Fi Thriller War Western				
1	Toy Story (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Toy%20Story%20(1995)	
		000011110000000000000000			
2	GoldenEye (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?GoldenEye%20(1995)	
		0011100000000000000000001000			
3	Four Rooms (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Four%20Rooms%20(1995)	
		0000000000000000000000001000			
4	Get Shorty (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Get%20Shorty%20(1995)	
		0010000100010000000000000000			
5	Copycat (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Copycat%20(1995)	
		0000000010010000000000001000			
6	Shanghai Triad (Yao a yao dao waipo qiao) (1995)	01-Jan-1995		http://us.imdb.com/Title?Yao+a+yao+yao+dao+waipo+qiao+(1995)	
		000000000000000010000000000000			
7	Twelve Monkeys (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Twelve%20Monkeys%20(1995)	
		00000000000010000000000010000			

1.4 Output Files

1.4.1 averageMovies.png



```
linux.cs.odu.edu - PuTTY
% clear ; python averageMovies.py
*****
Top 5 movies having the highest average ratings
*****
Movie Name                                Avg Rating
-----
They Made Me a Criminal                   5.0
Star Kid                                  5.0
Someone Else's America                    5.0
Santa with Muscles                        5.0
Saint of Fort Washington, The             5.0
% 
```

Figure 1: Top 5 movies with highest average ratings.

1.4.2 averageMoviesMore.txt

Movies with highest average ratings. The above output shows only top 5 but the below output shows all the movies with the same average ratings.

```
1 % clear ; python averageMovies.py
2 *****
3 Top 5 movies having the highest average ratings
4 *****
5 Movie Name                               Avg Rating
6 -----
7 They Made Me a Criminal                  5.0
8 Star Kid                                5.0
9 Someone Else 's America                  5.0
10 Santa with Muscles                      5.0
11 Saint of Fort Washington, The           5.0
12 Prefontaine                            5.0
13 Marlene Dietrich: Shadow and Light      5.0
14 Great Day in Harlem, A                  5.0
15 Entertaining Angels: The Dorothy Day Story 5.0
16 Aiqing wansui                           5.0
17 Pather Panchali                         4.6
18 Some Mother 's Son                      4.5
```

2 Question 2:

What 5 movies received the most ratings? Show the movies and the number of ratings sorted by number of ratings.

2.1 Approach

1. The approach is same as question 1.
2. The only difference is we count the number of ratings for each movie and display the corresponding movie name and number of ratings for each movie.
3. The output displays only top 5 movies and their number of ratings.
4. Figure 2 is the output for the program.

2.2 Source Code

2.2.1 mostRatingMovies.py

```
1 #!/usr/bin/env python
2
3 import sys
4
5 def main():
6     # Create a dictionary
7     movie_id_ratings = {}
8     movie_id_avg = {}
9     count = 1
10    # Read the input files.
11    readData = open('/home/bbokka/cs594/A8/u.data', 'r')
12    readItem = open('/home/bbokka/cs594/A8/u.item', 'r')
13    # reading the u.data file for item and ratings.
14    for line in readData:
15        split_input_line = line.strip().split('\t')
16        movie_id = split_input_line[1]
17        movie_rating = int(split_input_line[2])
18        try:
19            movie_id_ratings[movie_id].append(movie_rating)
20        except KeyError:
21            movie_id_ratings[movie_id] = list()
22            movie_id_ratings[movie_id].append(movie_rating)
23    readData.close()
24    # Calculating the average.
25    for key in movie_id_ratings:
26        length = len(movie_id_ratings[key])
27        movie_id_avg[key] = [int(length)]
28
29    # Reading the u.item file for movie name.
30    for each_line in readItem:
31        split_each_line = each_line.strip().split('|')
32        movie_item_id = split_each_line[0]
33        movie_name_split = split_each_line[1].split('(1')
34        movie_name = movie_name_split[0]
35
36        movie_id_avg[movie_item_id].append(movie_name)
37    readItem.close()
38    print '*' * 60
39    print "Top 5 movies received the most ratings"
40    print "*" * 60
41    print "Movie Name\t\t\t\t\t", "Number of ratings"
42    print "-" * 60
43    # sorting the movies from highest to lowest based on the average value.
44    for key, value in sorted(movie_id_avg.items(), key=lambda e: e[1], reverse=True):
45        if (count <= 5):
46            print '{:<50}{} ' .format(value[1], value[0])
47            count += 1
48
49 if __name__ == "__main__":
50     try:
51         main()
52     except KeyboardInterrupt:
53         sys.exit(1)
```

2.3 Input Files

2.3.1 u.data

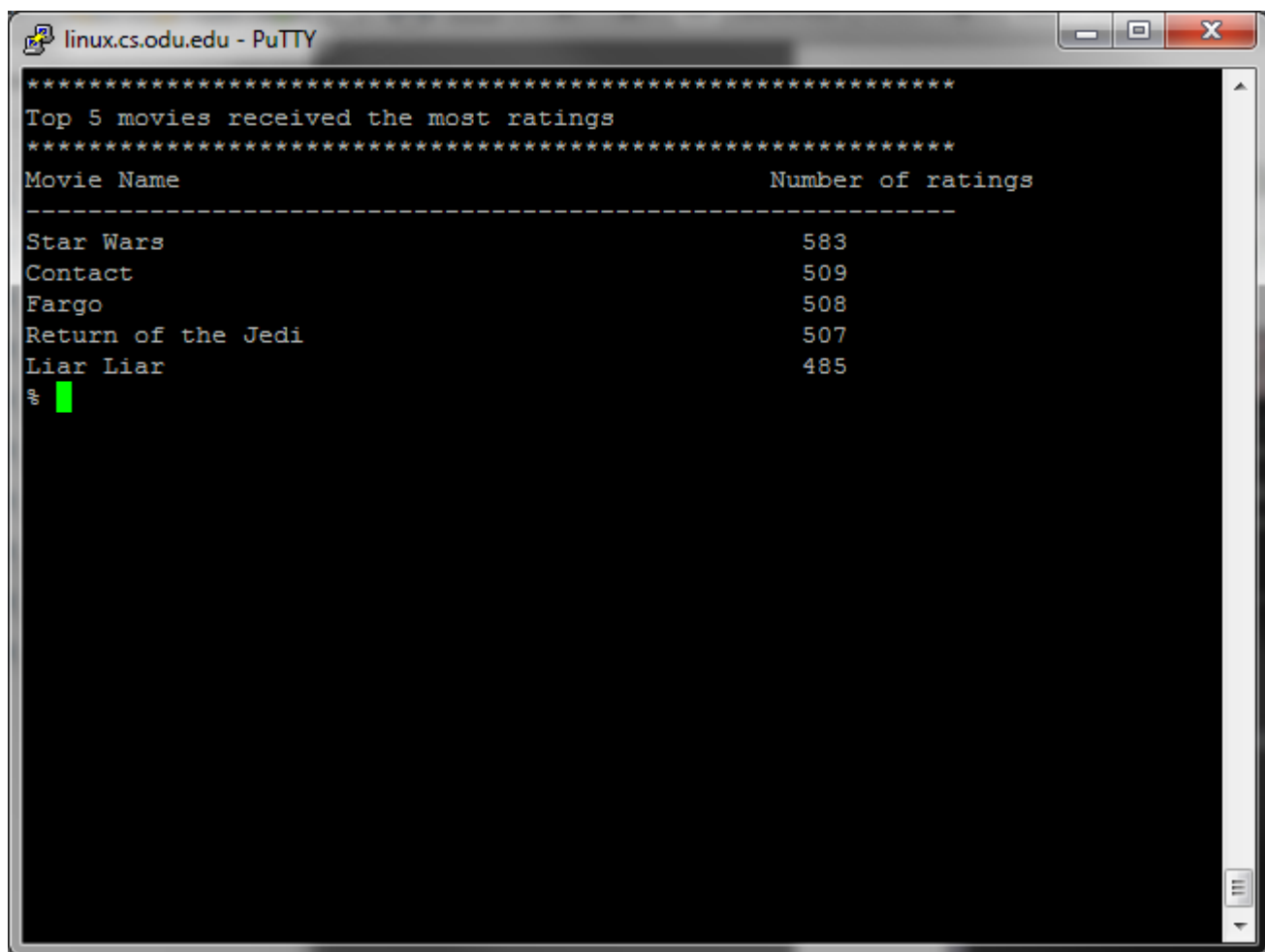
	userid	itemid	rating	timestamp
1				
2				
3	196	242	3	881250949
4	186	302	3	891717742
5	22	377	1	878887116
6	244	51	2	880606923
7	166	346	1	886397596
8	298	474	4	884182806
9	115	265	2	881171488

2.3.2 u.item

movie id	movie title	release date	video release date	IMDb URL	unknown
	Action Adventure Animation Children's Comedy Crime Documentary Drama Fantasy Film-Noir Horror Musical Mystery Romance Sci-Fi Thriller War Western				
1	Toy Story (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Toy%20Story%20(1995)	
		000011110000000000000000			
2	GoldenEye (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?GoldenEye%20(1995)	
		0011100000000000000000001000			
3	Four Rooms (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Four%20Rooms%20(1995)	
		0000000000000000000000001000			
4	Get Shorty (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Get%20Shorty%20(1995)	
		0010000100010000000000000000			
5	Copycat (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Copycat%20(1995)	
		0000000010010000000000001000			
6	Shanghai Triad (Yao a yao dao waipo qiao) (1995)	01-Jan-1995		http://us.imdb.com/Title?Yao+a+yao+yao+dao+waipo+qiao+(1995)	
		000000000000000010000000000000			
7	Twelve Monkeys (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Twelve%20Monkeys%20(1995)	
		00000000000010000000000010000			

2.4 Output Files

2.4.1 mostRatingMovies.png



```
linux.cs.odu.edu - PuTTY
*****
Top 5 movies received the most ratings
*****
Movie Name                                Number of ratings
-----
Star Wars                                583
Contact                                  509
Fargo                                    508
Return of the Jedi                       507
Liar Liar                                485
% 
```

Figure 2: Top 5 movies received the most ratings.

3 Question 3:

What 5 movies were rated the highest on average by women? Show the movies and their ratings sorted by ratings.

3.1 Approach

1. To solve this question the straight forward approach is to read u.user file get the user id, gender. Read u.data file get movie id, corresponding movie ratings. Read u.item file to get movie id and corresponding movie name.
2. From u.user file extract the user id and gender store them in a dictionary.
3. Now, Extract the movie id and ratings. For each user id build a dictionary that has all the ratings for the respective gender.
4. Now, for each movie id calculate the average a using `sum()` and `len()` functions.
5. For each movie id get the movie name and append it to the dictionary which has average.
6. Now, sort the dictionary and print the top 5 movies rated the highest on average by women.
7. Program `averageGender.py` is designed in such a way that it can get the highest averages by men and women just by changing the arguments to F or M while executing the program.
8. Figure 3 is the output showing top 5 movies rated the highest on average by women. The results of the program are limited to top 5 it can be modified to get more average ratings.

3.2 Source Code

3.2.1 averageMovies.py

```
1 #!/usr/bin/env python
2
3
4 # program to get the highest average ratings by men and women
5 import sys
6 import operator
7 def main():
8     # take the arguments from the command line
9     numOfArgs=len(sys.argv)
10    if numOfArgs<2 or numOfArgs>2:
11
12        print 'Usage: averageGender.py <F or M> '
13        print 'e.g. : averageGender.py F'
14        sys.exit(1)
15    gender = sys.argv[1]
16    # Create a dictionary
17    movie_userid_avg = {}
18    movie_id_ratings = {}
19    movie_userid_avg = {}
20    users = {}
21    count = 1
22
23    readData = open('/home/bbokka/cs594/A8/u.data', 'r')
24    readItem = open('/home/bbokka/cs594/A8/u.item', 'r')
25    readUser = open('/home/bbokka/cs594/A8/u.user', 'r')
26
27    # Reading u.user file to get the male or female
28    for line in readUser:
29        split_input_line = line.strip().split('|')
30        user_id_from_user= split_input_line[0]
31        user_gender = split_input_line[2]
32
33        if(gender.upper() == user_gender):
34            users[user_id_from_user] = [user_gender]
35
36    readUser.close()
37
38    # reading the u.data file for item and ratings.
39    for line in readData:
40        split_input_line = line.strip().split('\t')
41        user_id_from_data = split_input_line[0]
42        movie_id_from_data= split_input_line[1]
43        movie_rating = float(split_input_line[2])
44
45        if user_id_from_data in users:
46            try:
47                movie_id_ratings[movie_id_from_data].append(movie_rating)
48            except KeyError:
49                movie_id_ratings[movie_id_from_data] = list()
50                movie_id_ratings[movie_id_from_data].append(movie_rating)
51    readData.close()
52
53    # Calculating the average.
```

```

54     for key in movie_id_ratings:
55         avg = sum(movie_id_ratings[key])/len(movie_id_ratings[key])
56         movie_userid_avg[key] = [float(avg)]
57
58     # Reading the u.item file for movie name.
59     for each_line in readItem:
60         split_each_line = each_line.strip().split('|')
61         movie_item_id = split_each_line[0]
62         movie_name_split = split_each_line[1].split('(1')
63         movie_name = movie_name_split[0]
64
65         try:
66             movie_userid_avg[movie_item_id].append(movie_name)
67         except KeyError:
68             pass
69     readItem.close()
70     print "*" * 60
71     print "Top 5 movies rated the highest on average by women"
72     print "*" * 60
73     print "Movie Name\t\t\t\t\t", "Avg Rating"
74     print "-" * 60
75     # sorting the movies from highest to lowest based on the average value.
76     for key, value in sorted(movie_userid_avg.items(), key=lambda e: e[1], reverse=
True):
77         if (count <=20 ):
78             print '{:<50}{:<0.1f} ' .format(value[1], value[0])
79             count += 1
80
81 if __name__ == "__main__":
82     try:
83         main()
84     except KeyboardInterrupt:
85         sys.exit(1)

```


3.3 Input Files

3.3.1 u.data

	userid	itemid	rating	timestamp
1				
2				
3	196	242	3	881250949
4	186	302	3	891717742
5	22	377	1	878887116
6	244	51	2	880606923
7	166	346	1	886397596
8	298	474	4	884182806
9	115	265	2	881171488

3.3.2 u.item

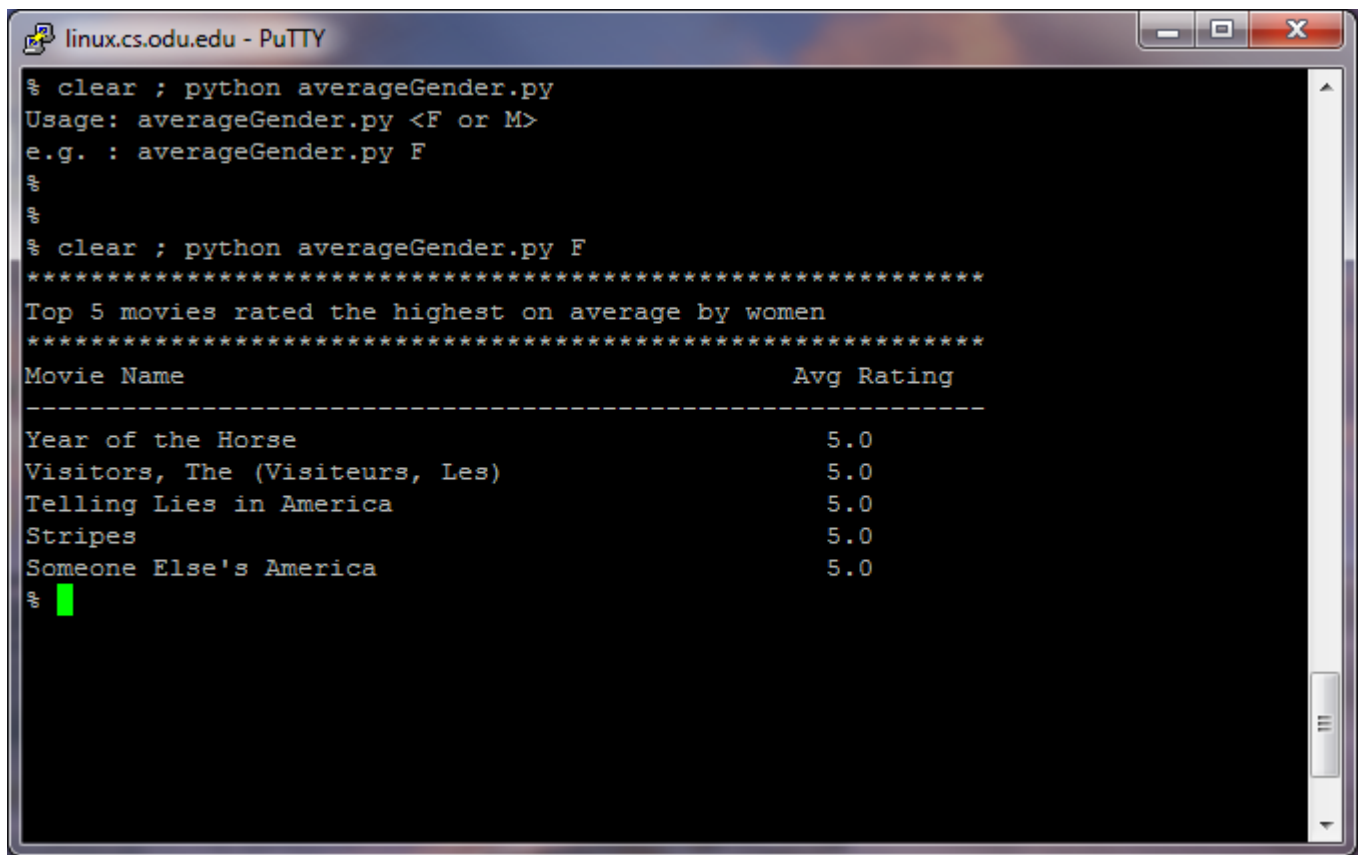
movie id	movie title	release date	video release date	IMDb URL	unknown
	Action Adventure Animation Children's Comedy Crime Documentary Drama Fantasy Film-Noir Horror Musical Mystery Romance Sci-Fi Thriller War Western				
1	Toy Story (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Toy%20Story%20(1995)	
2	GoldenEye (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?GoldenEye%20(1995)	
3	Four Rooms (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Four%20Rooms%20(1995)	
4	Get Shorty (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Get%20Shorty%20(1995)	
5	Copycat (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Copycat%20(1995)	
6	Shanghai Triad (Yao a yao dao waipo qiao) (1995)	01-Jan-1995		http://us.imdb.com/Title?Yao+a+yao+yao+dao+waipo+qiao+(1995)	
7	Twelve Monkeys (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Twelve%20Monkeys%20(1995)	

3.3.3 u.user

	userid	age	gender	occupation	zip_code
1	1	24	M	technician	85711
2	2	53	F	other	94043
3	3	23	M	writer	32067
4	4	24	M	technician	43537
5	5	33	F	other	15213
6	6	42	M	executive	98101
7	7	57	M	administrator	91344

3.4 Output Files

3.4.1 averageGenderF.png



The screenshot shows a PuTTY terminal window titled "linux.cs.odu.edu - PuTTY". The user has entered the command `% clear ; python averageGender.py`. The script displays usage instructions: `Usage: averageGender.py <F or M>` and `e.g. : averageGender.py F`. The user then enters `% clear ; python averageGender.py F`. The script outputs a separator line of asterisks, followed by the text "Top 5 movies rated the highest on average by women", another separator line, and a table with two columns: "Movie Name" and "Avg Rating". The table lists five movies, all with an average rating of 5.0. The terminal ends with a prompt `%` and a green cursor.

```
% clear ; python averageGender.py
Usage: averageGender.py <F or M>
e.g. : averageGender.py F
%
%
% clear ; python averageGender.py F
*****
Top 5 movies rated the highest on average by women
*****
Movie Name                               Avg Rating
-----
Year of the Horse                         5.0
Visitors, The (Visiteurs, Les)           5.0
Telling Lies in America                  5.0
Stripes                                  5.0
Someone Else's America                   5.0
% █
```

Figure 3: Top 5 movies rated the highest on average by women.

3.4.2 averageGenderFull.txt

Showing all the movies with same average ratings.

```
1 % clear ; python averageGender.py
2 Usage: averageGender.py <F or M>
3 e.g. : averageGender.py F
4 %
5 %
6 %
7 % clear ; python averageGender.py F
8 *****
9 Top 5 movies rated the highest on average by women
10 *****
11 Movie Name                               Avg Rating
12 -----
13 Year of the Horse                        5.0
14 Visitors , The (Visiteurs , Les)        5.0
15 Telling Lies in America                 5.0
16 Stripes                                 5.0
17 Someone Else 's America                 5.0
18 Prefontaine                             5.0
19 Mina Tannenbaum                          5.0
20 Maya Lin: A Strong Clear Vision          5.0
21 Foreign Correspondent                    5.0
22 Faster Pussycat! Kill! Kill!            5.0
23 Everest                                 5.0
24 Schindler 's List                        4.6
25 Close Shave , A                         4.6
26 Shawshank Redemption , The              4.6
27 Wallace & Gromit: The Best of Aardman Animation 4.5
28 Shall We Dance?                         4.5
29 Women, The                              4.5
30 Some Mother 's Son                      4.5
31 Friday                                  4.5
32 Band Wagon, The                         4.5
```

4 Question 4:

What 5 movies were rated the highest on average by men? Show the movies and their ratings sorted by ratings.

4.1 Approach

1. The approach is same as question 3.
2. The only difference is changing the command line argument which can be seen in the output produced.
3. The output displays only top 5 movies rated the highest on average by men.
4. Figure 4 is the output for the program.

4.2 Source Code

4.2.1 averageMovies.py

```
1 #!/usr/bin/env python
2
3
4 # program to get the highest average ratings by men and women
5 import sys
6 import operator
7 def main():
8     # take the arguments from the command line
9     numOfArgs=len(sys.argv)
10    if numOfArgs<2 or numOfArgs>2:
11
12        print 'Usage: averageGender.py <F or M> '
13        print 'e.g. : averageGender.py F'
14        sys.exit(1)
15    gender = sys.argv[1]
16    # Create a dictionary
17    movie_userid_avg = {}
18    movie_id_ratings = {}
19    movie_userid_avg = {}
20    users = {}
21    count = 1
22
23    readData = open('/home/bbokka/cs594/A8/u.data', 'r')
24    readItem = open('/home/bbokka/cs594/A8/u.item', 'r')
25    readUser = open('/home/bbokka/cs594/A8/u.user', 'r')
26
27    # Reading u.user file to get the male or female
28    for line in readUser:
29        split_input_line = line.strip().split('|')
30        user_id_from_user= split_input_line[0]
31        user_gender = split_input_line[2]
32
33        if(gender.upper() == user_gender):
34            users[user_id_from_user] = [user_gender]
35
36    readUser.close()
37
38    # reading the u.data file for item and ratings.
39    for line in readData:
40        split_input_line = line.strip().split('\t')
41        user_id_from_data = split_input_line[0]
42        movie_id_from_data= split_input_line[1]
43        movie_rating = float(split_input_line[2])
44
45        if user_id_from_data in users:
46            try:
47                movie_id_ratings[movie_id_from_data].append(movie_rating)
48            except KeyError:
49                movie_id_ratings[movie_id_from_data] = list()
50                movie_id_ratings[movie_id_from_data].append(movie_rating)
51    readData.close()
52
53    # Calculating the average.
```

```

54     for key in movie_id_ratings:
55         avg = sum(movie_id_ratings[key])/len(movie_id_ratings[key])
56         movie_userid_avg[key] = [float(avg)]
57
58     # Reading the u.item file for movie name.
59     for each_line in readItem:
60         split_each_line = each_line.strip().split('|')
61         movie_item_id = split_each_line[0]
62         movie_name_split = split_each_line[1].split('(1')
63         movie_name = movie_name_split[0]
64
65         try:
66             movie_userid_avg[movie_item_id].append(movie_name)
67         except KeyError:
68             pass
69     readItem.close()
70     print "*" * 60
71     print "Top 5 movies rated the highest on average by women"
72     print "*" * 60
73     print "Movie Name\t\t\t\t\t", "Avg Rating"
74     print "-" * 60
75     # sorting the movies from highest to lowest based on the average value.
76     for key, value in sorted(movie_userid_avg.items(), key=lambda e: e[1], reverse=
77 True):
78         if (count <=20 ):
79             print '{:<50}{:<0.1f} '.format(value[1], value[0])
80             count += 1
81
82 if __name__ == "__main__":
83     try:
84         main()
85     except KeyboardInterrupt:
86         sys.exit(1)

```

4.3 Input Files

4.3.1 u.data

	userid	itemid	rating	timestamp
1				
2				
3	196	242	3	881250949
4	186	302	3	891717742
5	22	377	1	878887116
6	244	51	2	880606923
7	166	346	1	886397596
8	298	474	4	884182806
9	115	265	2	881171488

4.3.2 u.item

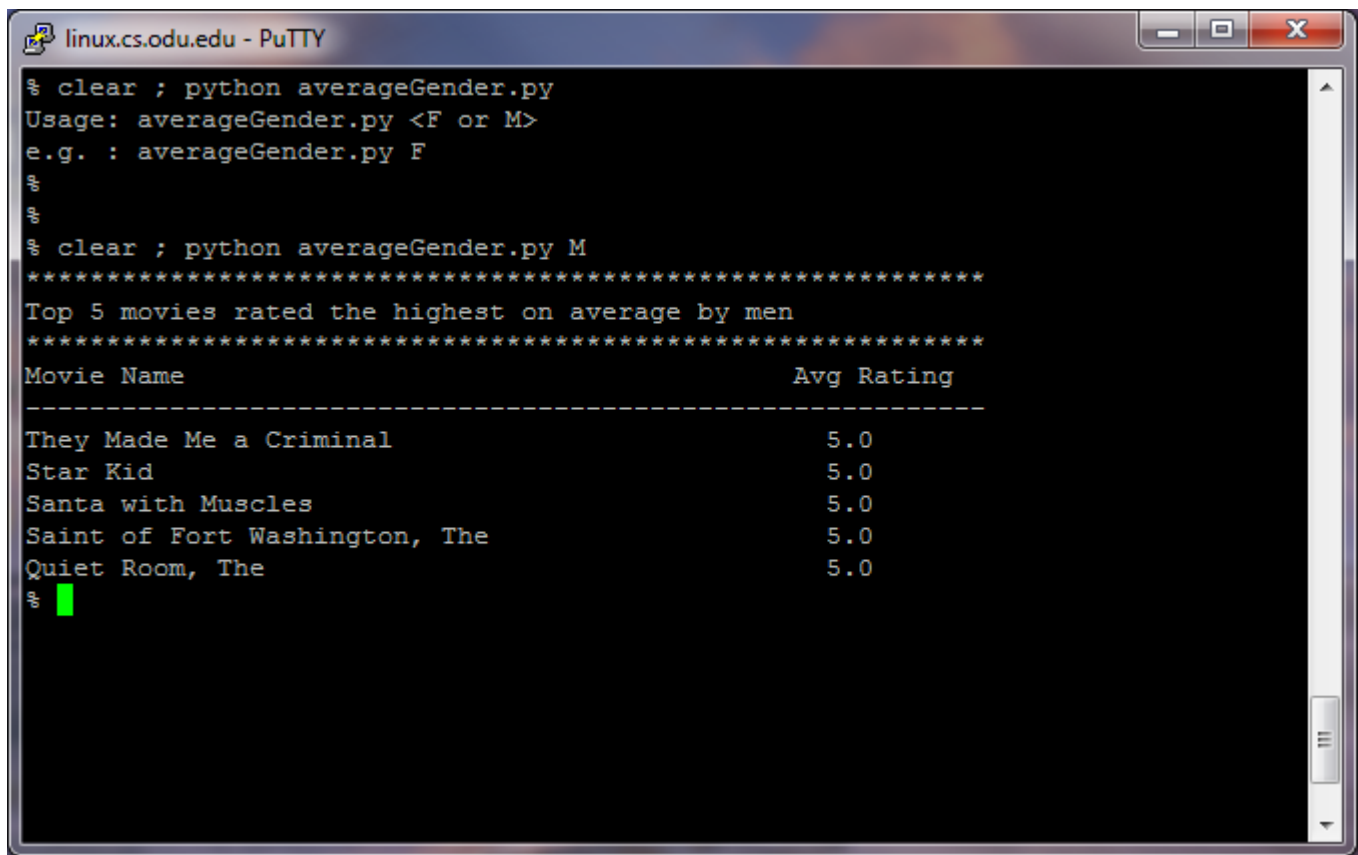
movie id	movie title	release date	video release date	IMDb URL	unknown
	Action Adventure Animation Children's Comedy Crime Documentary Drama Fantasy Film-Noir Horror Musical Mystery Romance Sci-Fi Thriller War Western				
1	Toy Story (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Toy%20Story%20(1995)	
		000011110000000000000000			
2	GoldenEye (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?GoldenEye%20(1995)	
		0011100000000000000000001000			
3	Four Rooms (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Four%20Rooms%20(1995)	
		0000000000000000000000001000			
4	Get Shorty (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Get%20Shorty%20(1995)	
		0010000100010000000000000000			
5	Copycat (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Copycat%20(1995)	
		0000000010010000000000001000			
6	Shanghai Triad (Yao a yao dao waipo qiao) (1995)	01-Jan-1995		http://us.imdb.com/Title?Yao+a+yao+yao+dao+waipo+qiao+(1995)	
		000000000000000010000000000000			
7	Twelve Monkeys (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Twelve%20Monkeys%20(1995)	
		00000000000010000000000010000			

4.3.3 u.user

	userid	age	gender	occupation	zip_code
1	1	24	M	technician	85711
2	2	53	F	other	94043
3	3	23	M	writer	32067
4	4	24	M	technician	43537
5	5	33	F	other	15213
6	6	42	M	executive	98101
7	7	57	M	administrator	91344

4.4 Output Files

4.4.1 averageGenderM.png



```
linux.cs.odu.edu - PuTTY
% clear ; python averageGender.py
Usage: averageGender.py <F or M>
e.g. : averageGender.py F
%
%
% clear ; python averageGender.py M
*****
Top 5 movies rated the highest on average by men
*****
Movie Name                               Avg Rating
-----
They Made Me a Criminal                  5.0
Star Kid                                 5.0
Santa with Muscles                       5.0
Saint of Fort Washington, The            5.0
Quiet Room, The                          5.0
% █
```

Figure 4: Top 5 movies rated the highest on average by men.

4.4.2 averageGenderFull.txt

Showing all the movies with same average ratings.

```
1 % clear ; python averageGender.py
2 Usage: averageGender.py <F or M>
3 e.g. : averageGender.py F
4 %
5 %
6 % clear ; python averageGender.py M
7 *****
8 Top 5 movies rated the highest on average by women
9 *****
10 Movie Name                               Avg Rating
11 -----
12 They Made Me a Criminal                 5.0
13 Star Kid                               5.0
14 Santa with Muscles                     5.0
15 Saint of Fort Washington, The          5.0
16 Quiet Room, The                        5.0
17 Prefontaine                           5.0
18 Marlene Dietrich: Shadow and Light     5.0
19 Love Serenade                          5.0
20 Little City                            5.0
21 Letter From Death Row, A               5.0
22 Leading Man, The                       5.0
23 Hugo Pool                              5.0
24 Great Day in Harlem, A                 5.0
25 Entertaining Angels: The Dorothy Day Story 5.0
26 Delta of Venus                         5.0
27 Aiqing wansui                          5.0
28 Two or Three Things I Know About Her   4.7
29 Little Princess, The                   4.7
30 Pather Panchali                        4.6
31 Sliding Doors                          4.5
```

5 Question 5:

What movie received ratings most like Top Gun? Which movie received ratings that were least like Top Gun (negative correlation)?

5.1 Approach

1. To solve question 5 we need recommendation.py.
2. To get the desired output I modified calculateSimilarItems() and loadMovieLens() functions.
3. The file ratingsLikeTopGun.txt shows the output for movies having ratings most like Top Gun.
4. The file ratingsLeastTopGun.txt shows the output for movies having ratings least like Top Gun.
5. To see more results run ratingsTopGun.py.

5.2 Source Code

5.2.1 recommendations.py

```
1 def calculateSimilarItems(prefs,n=10):
2     # Create a dictionary of items showing which other items they
3     # are most similar to.
4     result={}
5     item1 = 'Top Gun (1986)'
6     # Invert the preference matrix to be item-centric
7     itemPrefs=transformPrefs(prefs)
8     c=0
9     for item in itemPrefs:
10        # Status updates for large datasets
11        # c+=1
12        # if c%100==0: print "%d / %d" % (c,len(itemPrefs))
13        # Find the most similar items to this one
14        scores=topMatches(itemPrefs,item1,n=n,similarity=sim_pearson)
15        result[item]=scores
16    return result
17
18
19 def loadMovieLens():
20     # Get movie titles
21     readData = open('/home/bbokka/cs594/A8/u.data', 'r')
22     readItem = open('/home/bbokka/cs594/A8/u.item', 'r')
23     movies={}
24     for line in readItem:
25         (id,title)=line.split('|')[0:2]
26         movies[id]=title
27     # Load data
28     prefs={}
29     for line in readData:
30         (user,movieid,rating,ts)=line.split('\t')
31         prefs.setdefault(user,{})
32         prefs[user][movies[movieid]]=float(rating)
33
34     return prefs
35
36 # To get the negative correlation line 42 is commented.
37 def topMatches(prefs,person,n=5,similarity=sim_pearson):
38     scores=[(similarity(prefs,person,other),other)
39             for other in prefs if other!=person]
40     scores.sort()
41     #scores.reverse()
42     return scores[0:n]
```

5.2.2 ratingsTopGun.py

```
1 #!/usr/bin/env python
2
3 import sys
4 import recommendations
5
6 def main():
7     count = 'Top Gun (1986)'
8     resultsOfloadMovieLens = recommendations.loadMovieLens()
9     resultsOfcalculateSimilarItems = recommendations.calculateSimilarItems(
10     resultsOfloadMovieLens,n=80)
11     print "*" * 60
12     print "Movies received ratings most like or least like Top Gun"
13     print "*" * 60
14     print "value\t\t\t\t\t", "Movie Name"
15     print "-" * 60
16
17     for key, value in sorted(resultsOfcalculateSimilarItems.items(), key=lambda e: e
18     [1], reverse=True):
19         variable = key
20         if (count == variable):
21             for value , movie in value:
22                 print value , movie
23
24 if __name__ == "__main__":
25     try:
26         main()
27     except KeyboardInterrupt:
28         sys.exit(1)
```

5.3 Input Files

5.3.1 u.data

	userid	itemid	rating	timestamp
1				
2				
3	196	242	3	881250949
4	186	302	3	891717742
5	22	377	1	878887116
6	244	51	2	880606923
7	166	346	1	886397596
8	298	474	4	884182806
9	115	265	2	881171488

5.3.2 u.item

movie id	movie title	release date	video release date	IMDb URL	unknown
	Action Adventure Animation Children's Comedy Crime Documentary Drama Fantasy Film-Noir Horror Musical Mystery Romance Sci-Fi Thriller War Western				
1	Toy Story (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Toy%20Story%20(1995)	
		000011110000000000000000			
2	GoldenEye (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?GoldenEye%20(1995)	
		0011100000000000000000001000			
3	Four Rooms (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Four%20Rooms%20(1995)	
		0000000000000000000000001000			
4	Get Shorty (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Get%20Shorty%20(1995)	
		0010000100010000000000000000			
5	Copycat (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Copycat%20(1995)	
		0000000010010000000000001000			
6	Shanghai Triad (Yao a yao dao waipo qiao) (1995)	01-Jan-1995		http://us.imdb.com/Title?Yao+a+yao+yao+dao+waipo+qiao+(1995)	
		000000000000000010000000000000			
7	Twelve Monkeys (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Twelve%20Monkeys%20(1995)	
		00000000000010000000000010000			

5.4 Output Files

5.4.1 ratingsLikeTopGun.txt

```
1 % clear ; python ratingsTopGun.py
2 *****
3 Movies received ratings most like Top Gun
4 *****
5 value                               Movie Name
6 -----
7 1.0 Shiloh (1997)
8 1.0 King of the Hill (1993)
9 1.0 Bhaji on the Beach (1993)
10 1.0 Wild America (1997)
11 1.0 Wedding Gift , The (1994)
12 1.0 Underground (1995)
13 1.0 Two or Three Things I Know About Her (1966)
14 1.0 Two Bits (1995)
15 1.0 Total Eclipse (1995)
16 1.0 The Innocent (1994)
17 1.0 That Old Feeling (1997)
18 1.0 Stars Fell on Henrietta , The (1995)
19 1.0 Stalker (1979)
20 1.0 Spirits of the Dead (Tre passi nel delirio) (1968)
21 1.0 Show, The (1995)
22 1.0 Shooter, The (1995)
23 1.0 Selena (1997)
24 1.0 Schizopolis (1996)
25 1.0 Scarlet Letter, The (1926)
26 1.0 Run of the Country, The (1995)
27 1.0 Ponette (1996)
28 1.0 Perfect Candidate, A (1996)
29 1.0 Outlaw, The (1943)
30 1.0 Old Lady Who Walked in the Sea, The (Vieille qui marchait dans la mer, La) (1991)
31 1.0 Nothing to Lose (1994)
32 1.0 New Jersey Drive (1995)
33 1.0 Mr. Jones (1993)
34 1.0 Metisse ( Caf au Lait) (1993)
35 1.0 Maybe, Maybe Not (Bewegte Mann, Der) (1994)
36 1.0 Manny & Lo (1996)
37 1.0 Man of the Year (1995)
38 1.0 Love Serenade (1996)
39 1.0 Last Time I Saw Paris, The (1954)
40 1.0 Killer (Bulletproof Heart) (1994)
41 1.0 Jerky Boys, The (1994)
42 1.0 I Like It Like That (1994)
43 1.0 Horse Whisperer, The (1998)
44 1.0 Hear My Song (1991)
45 1.0 Grosse Fatigue (1994)
46 1.0 Gone Fishin ' (1997)
47 1.0 Glass Shield, The (1994)
48 1.0 Germinal (1993)
49 1.0 Gabbeh (1996)
50 1.0 Four Days in September (1997)
51 1.0 Flower of My Secret, The (Flor de mi secreto, La) (1995)
52 1.0 Fausto (1993)
53 1.0 Even Cowgirls Get the Blues (1993)
```

54 1.0 Enfer , L' (1994)
55 1.0 Dream With the Fishes (1997)
56 1.0 Dream Man (1995)
57 1.0 Dangerous Ground (1997)
58 1.0 Collectionneuse , La (1967)
59 1.0 Clean Slate (Coup de Torchon) (1981)
60 1.0 Calendar Girl (1993)
61 1.0 Blood For Dracula (Andy Warhol's Dracula) (1974)
62 1.0 Bliss (1997)
63 1.0 Best Men (1997)
64 1.0 American Dream (1990)
65 1.0 Albino Alligator (1996)
66 1.0 8 Seconds (1994)
67 1.0 Aparajito (1956)
68 0.995870594886 Scarlet Letter , The (1995)
69 0.986440050416 Kundun (1997)
70 0.981980506062 Traveller (1997)

5.4.2 ratingsLeastTopGun.txt

```
1 % clear ; python ratingsTopGun.py
2 *****
3 Movies received ratings most Least like Top Gun
4 *****
5 value                               Movie Name
6 -----
7 -1.0 Babysitter , The (1995)
8 -1.0 Telling Lies in America (1997)
9 -1.0 Bad Moon (1996)
10 -1.0 Beat the Devil (1954)
11 -1.0 Bewegte Mann, Der (1994)
12 -1.0 Bitter Sugar (Azucar Amargo) (1996)
13 -1.0 Broken English (1996)
14 -1.0 Caro Diario (Dear Diary) (1994)
15 -1.0 Carpool (1996)
16 -1.0 Carried Away (1996)
17 -1.0 Everest (1998)
18 -1.0 Frisk (1995)
19 -1.0 Heidi Fleiss: Hollywood Madam (1995)
20 -1.0 Joy Luck Club, The (1993)
21 -1.0 Lamerica (1994)
22 -1.0 Loch Ness (1995)
23 -1.0 Love and Death on Long Island (1997)
24 -1.0 Lover's Knot (1996)
25 -1.0 Meet Wally Sparks (1997)
26 -1.0 Midnight Dancers (Sibak) (1994)
27 -1.0 Naked in New York (1994)
28 -1.0 Nico Icon (1995)
29 -1.0 Nil By Mouth (1997)
30 -1.0 Romper Stomper (1992)
31 -1.0 Roseanna's Grave (For Roseanna) (1997)
32 -1.0 Safe Passage (1994)
33 -1.0 Switchback (1997)
34 -1.0 Tetsuo II: Body Hammer (1992)
35 -1.0 Two Much (1996)
36 -1.0 World of Apu, The (Apu Sansar) (1959)
37 -1.0 Year of the Horse (1997)
38 -0.946729262406 Alphaville (1965)
39 -0.944911182523 Denise Calls Up (1995)
40 -0.944911182523 Penny Serenade (1941)
41 -0.894427191 Awfully Big Adventure, An (1995)
42 -0.888431093553 Paradise Lost: The Child Murders at Robin Hood Hills (1996)
43 -0.880704845928 In the Realm of the Senses (Ai no corrida) (1976)
44 -0.866025403784 Chasers (1994)
45 -0.866025403784 Faces (1968)
46 -0.866025403784 Big Squeeze, The (1996)
47 -0.866025403784 Female Perversions (1996)
48 -0.816496580928 Theodore Rex (1995)
49 -0.777777777778 Half Baked (1998)
50 -0.777777777778 Wild Reeds (1994)
51 -0.770317441193 Cemetery Man (Dellamorte Dellamore) (1994)
52 -0.760885910253 Gang Related (1997)
53 -0.753266379297 Young Poisoner's Handbook, The (1995)
54 -0.751809411556 Ruby in Paradise (1993)
55 -0.734031784566 Funeral, The (1996)
```


56 -0.73029674334 Mark of Zorro , The (1940)
57 -0.73029674334 One Night Stand (1997)
58 -0.707106781187 Angel and the Badman (1947)
59 -0.688247201612 Switchblade Sisters (1975)
60 -0.683302297676 Apostle , The (1997)

6 Question 6:

Which 5 raters rated the most films? Show the raters' IDs and the number of films each rated.

6.1 Approach

1. The approach for this question is modify question 1 to read U.data file, get user id and for the corresponding user id store all the ratings .
2. Now, we count the number of ratings for each user by using the len() function.
3. The output displays only top 5 users and their number of ratings.
4. Figure 5 is the output for the program.

6.2 Source Code

6.2.1 ratersMostRatedMovies.py

```
1#!/usr/bin/env python
2
3import sys
4
5def main():
6    # Create a dictionary
7    movie_userid_ratings = {}
8    movie_userid_len     = {}
9    count                = 1
10   # Read the input files.
11   readData = open('/home/bbokka/cs594/A8/u.data', 'r')
12   # reading the u.data file for item and ratings.
13   for line in readData:
14       split_input_line = line.strip().split('\t')
15       user_id_from_data = split_input_line[0]
16       movie_rating      = float(split_input_line[2])
17
18       try:
19           movie_userid_ratings[user_id_from_data].append(movie_rating)
20       except KeyError:
21           movie_userid_ratings[user_id_from_data] = list()
22           movie_userid_ratings[user_id_from_data].append(movie_rating)
23   readData.close()
24   # Calculating the number of movies each user rated.
25   for key in movie_userid_ratings:
26       length = len(movie_userid_ratings[key])
27       movie_userid_len[key] = int(length)
28
29
30   print '*' * 60
31   print "Top 5 raters rated movies."
32   print "*" * 60
33   print "User id\t\t","Number of movies rated"
34   print "-" * 60
35   # sorting the movies from highest to lowest based on the number of
36   # ratings for each movie.
37   for key, value in sorted(movie_userid_len.items(), key=lambda e: e[1], reverse=
True):
38       if (count <= 5):
39           print '{:<20}{} ' .format(key, value)
40           count += 1
41
42 if __name__ == "__main__":
43     try:
44         main()
45     except KeyboardInterrupt:
46         sys.exit(1)
```

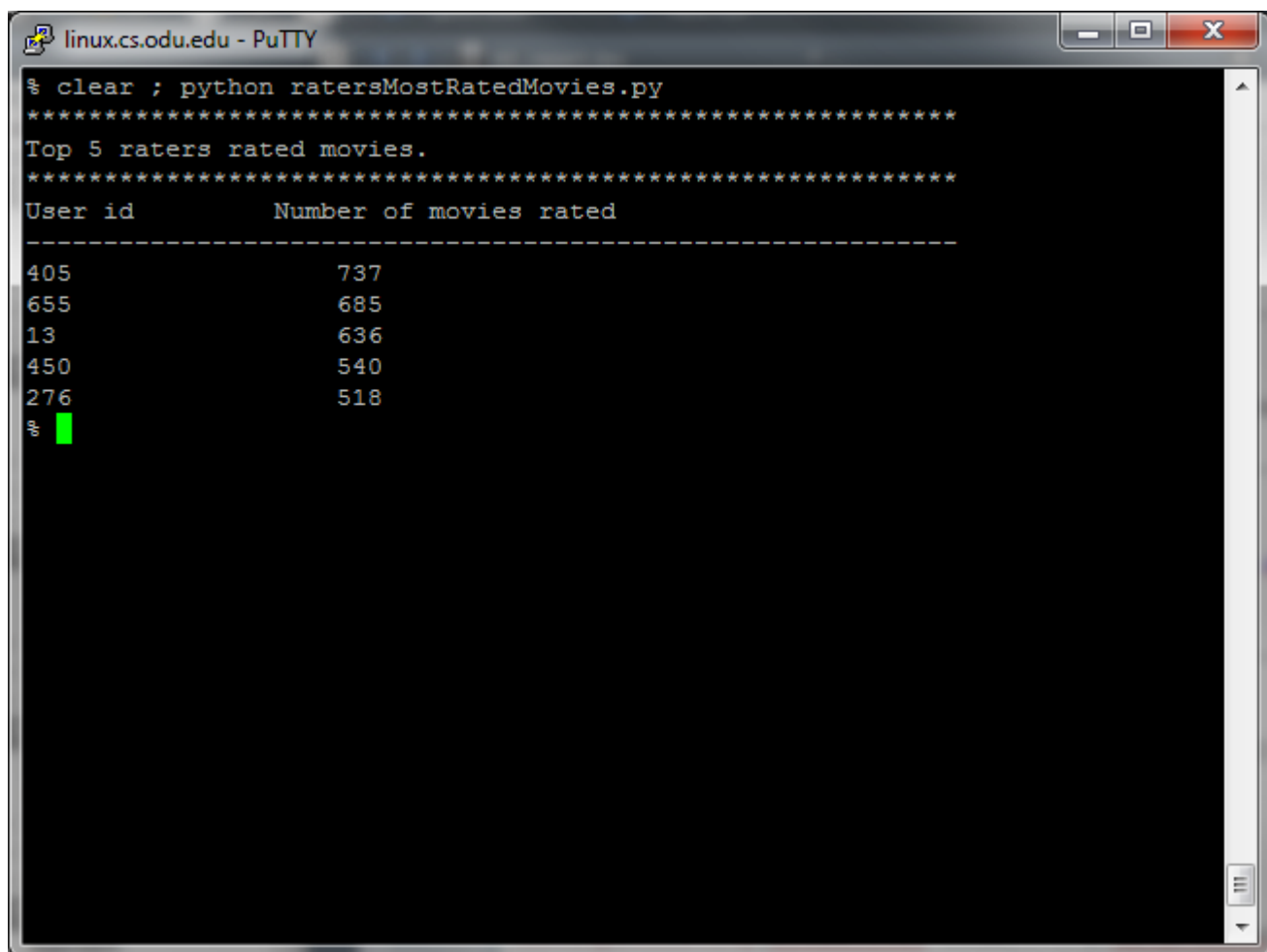
6.3 Input Files

6.3.1 u.data

```
1 userid itemid rating timestamp
2
3 196      242      3    881250949
4 186      302      3    891717742
5 22       377      1    878887116
6 244      51       2    880606923
7 166      346      1    886397596
8 298      474      4    884182806
9 115      265      2    881171488
```

6.4 Output Files

6.4.1 ratersMostRatedMovies.png



```
linux.cs.odu.edu - PuTTY
% clear ; python ratersMostRatedMovies.py
*****
Top 5 raters rated movies.
*****
User id      Number of movies rated
-----
405          737
655          685
13           636
450          540
276          518
% 
```

Figure 5: Top 5 raters rated movies.

7 Question 7:

Which 5 raters most agreed with each other? Show the raters' IDs and Pearson's r , sorted by r .

7.1 Approach

I couldn't complete this question.

7.2 Description of

7.3 Source Code

7.4 Output Files

8 Question 8:

Which 5 raters most disagreed with each other (negative correlation)? Show the raters' IDs and Pearson's r , sorted by r .

8.1 Approach

I couldn't complete this question.

8.2 Description of

8.3 Source Code

8.4 Output Files

9 Question 9:

What movie was rated highest on average by men over 40? By men under 40?

9.1 Approach

1. To solve this question, the straight forward approach is to read u.user file, get the user id, gender, age of each user. Read u.data file, get movie id, corresponding movie ratings. Read u.item file to get movie id and corresponding movie name.
2. Program averageOverUnderGender.py is designed in such a way that it can get the highest averages by men and women over and under 40 just by changing the arguments to F or M and “over” or “under” while executing the program.
3. From u.user file extract the user id, gender, age, based on arguments given from command line either men over 40 or men under 40 or women over 40 or women under 40 is stored into dictionary.
4. Now, Extract the movie id and ratings. For each user id build a dictionary that has all the ratings for the respective gender with over or under condition.
5. Now, for each movie id calculate the average a using sum() and len() functions.
6. For each movie id get the movie name and append it to the dictionary which has average.
7. Now, sort the dictionary and get the movies which are rated highest on average by men or women over 40 or under 40.
8. Figure 6 is the output showing how to execute the program and displays the top 5 movies which are rated highest on average by men “over” 40.
9. Figure 7 is the output showing how to execute the program and displays the top 5 movies which are rated highest on average by men “under” 40.
10. The file averageOverM.txt and averageUnderM.txt displays movies which have the same average ratings.

9.2 Source Code

9.2.1 averageOverUnderGender.py

```
1 #!/usr/bin/env python
2
3
4 # program to get the highest average ratings by men under and over 40.
5 #
6 import sys
7 import operator
8 def main():
9     # Take the arguments from the command line
10    numOfArgs=len(sys.argv)
11    if numOfArgs<3 or numOfArgs>3:
12
13        print 'Usage: averageOverUnderGender.py <F or M> <under or over>'
14        print 'e.g. : averageOverUnderGender.py F over '
15        sys.exit(1)
16    gender          = sys.argv[1]
17    gender_under_over= sys.argv[2]
18
19    # Create a dictionary
20    movie_userid_avg = {}
21    movie_id_ratings = {}
22    movie_userid_avg = {}
23    users            = {}
24    count            = 1
25
26    # Read the files
27    readData = open('/home/bbokka/cs594/A8/u.data', 'r')
28    readItem = open('/home/bbokka/cs594/A8/u.item', 'r')
29    readUser = open('/home/bbokka/cs594/A8/u.user', 'r')
30
31    # Reading u.user file to get the male or female
32    for line in readUser:
33        split_input_line = line.strip().split('|')
34        user_id_from_user= split_input_line[0]
35        user_age          = int(split_input_line[1])
36        user_gender       = split_input_line[2]
37
38        if gender_under_over == 'over' :
39            if(gender.upper() == user_gender and user_age > 40):
40                users[user_id_from_user] = [user_age]
41
42        if gender_under_over == 'under' :
43            if(gender.upper() == user_gender and user_age < 40):
44                users[user_id_from_user] = [user_age]
45
46    readUser.close()
47
48    # reading the u.data file for item and ratings.
49    for line in readData:
50        split_input_line = line.strip().split('\t')
51        user_id_from_data = split_input_line[0]
52        movie_id_from_data= split_input_line[1]
53        movie_rating      = float(split_input_line[2])
```



```

54
55     if user_id_from_data in users:
56         try:
57             movie_id_ratings[movie_id_from_data].append(movie_rating)
58         except KeyError:
59             movie_id_ratings[movie_id_from_data] = list()
60             movie_id_ratings[movie_id_from_data].append(movie_rating)
61 readData.close()
62
63 # Calculating the average.
64 for key in movie_id_ratings:
65     avg = sum(movie_id_ratings[key])/len(movie_id_ratings[key])
66     movie_userid_avg[key] = [float(avg)]
67
68 # Reading the u.item file for movie name.
69 for each_line in readItem:
70     split_each_line = each_line.strip().split('|')
71
72     movie_item_id = split_each_line[0]
73     movie_name_split = split_each_line[1].split('(1')
74     movie_name = movie_name_split[0]
75
76     try:
77         movie_userid_avg[movie_item_id].append(movie_name)
78     except KeyError:
79         pass
80 readItem.close()
81 print '*' * 55
82 print "Movies rated highest on average by men"
83 print "*" * 55
84 print "Movie Name\t\t\t\t\t","Avg Rating"
85 print "-" * 55
86 # sorting the movies from highest to lowest based on the average value.
87 for key, value in sorted(movie_userid_avg.items(), key=lambda e: e[1], reverse=
True):
88     if (count <=30 ):
89         print '{:<50}{:<0.1f} ' .format(value[1], value[0])
90         count += 1
91
92 if __name__ == "__main__":
93     try:
94         main()
95     except KeyboardInterrupt:
96         sys.exit(1)

```

9.3 Input Files

9.3.1 u.data

	userid	itemid	rating	timestamp
1				
2				
3	196	242	3	881250949
4	186	302	3	891717742
5	22	377	1	878887116
6	244	51	2	880606923
7	166	346	1	886397596
8	298	474	4	884182806
9	115	265	2	881171488

9.3.2 u.item

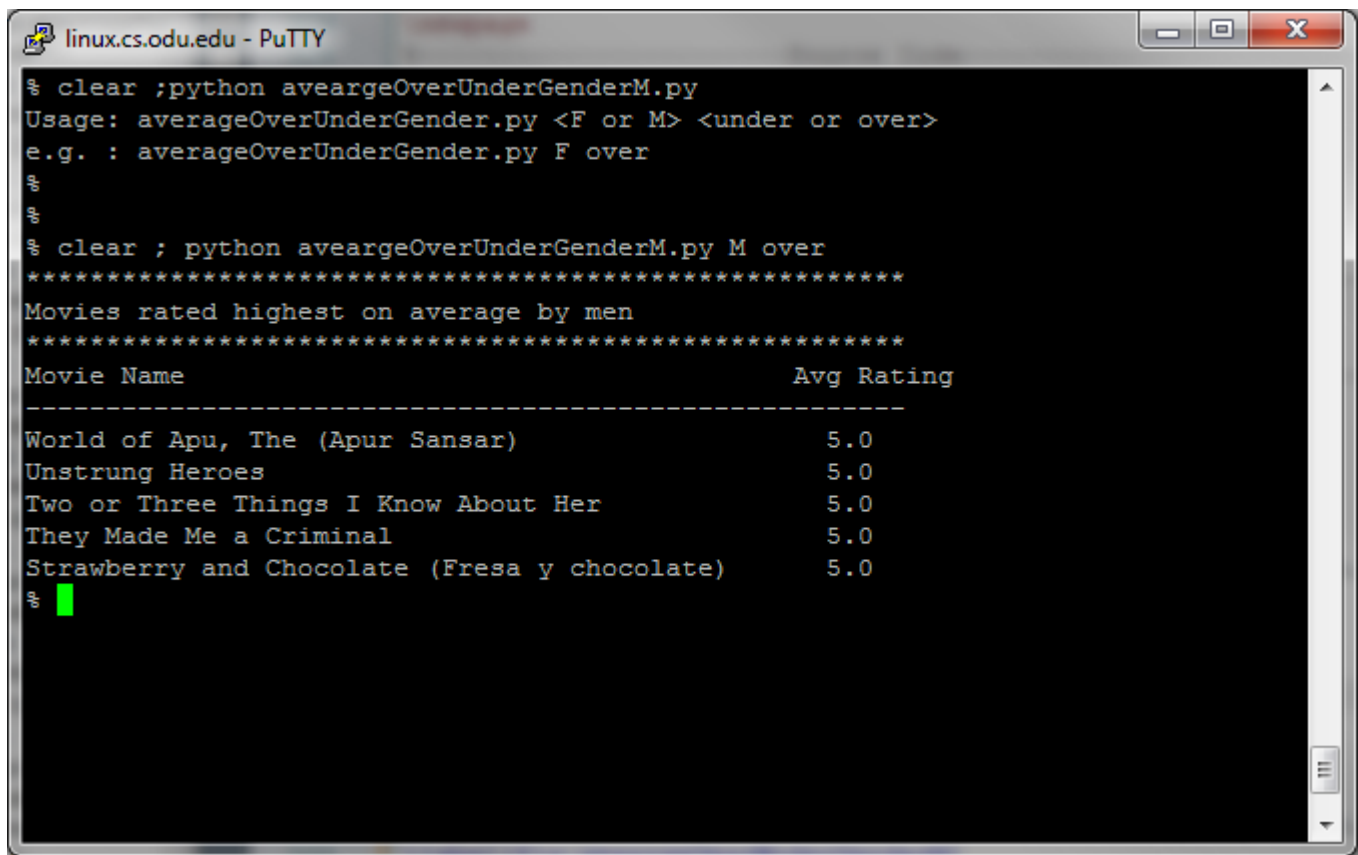
movie id	movie title	release date	video release date	IMDb URL	unknown
	Action Adventure Animation Children's Comedy Crime Documentary Drama Fantasy Film-Noir Horror Musical Mystery Romance Sci-Fi Thriller War Western				
1	Toy Story (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Toy%20Story%20(1995)	
2	GoldenEye (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?GoldenEye%20(1995)	
3	Four Rooms (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Four%20Rooms%20(1995)	
4	Get Shorty (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Get%20Shorty%20(1995)	
5	Copycat (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Copycat%20(1995)	
6	Shanghai Triad (Yao a yao yao dao waipo qiao) (1995)	01-Jan-1995		http://us.imdb.com/Title?Yao+a+yao+yao+dao+waipo+qiao+(1995)	
7	Twelve Monkeys (1995)	01-Jan-1995		http://us.imdb.com/M/title-exact?Twelve%20Monkeys%20(1995)	

9.3.3 u.user

	userid	age	gender	occupation	zip_code
1	1	24	M	technician	85711
2	2	53	F	other	94043
3	3	23	M	writer	32067
4	4	24	M	technician	43537
5	5	33	F	other	15213
6	6	42	M	executive	98101
7	7	57	M	administrator	91344

9.4 Output Files

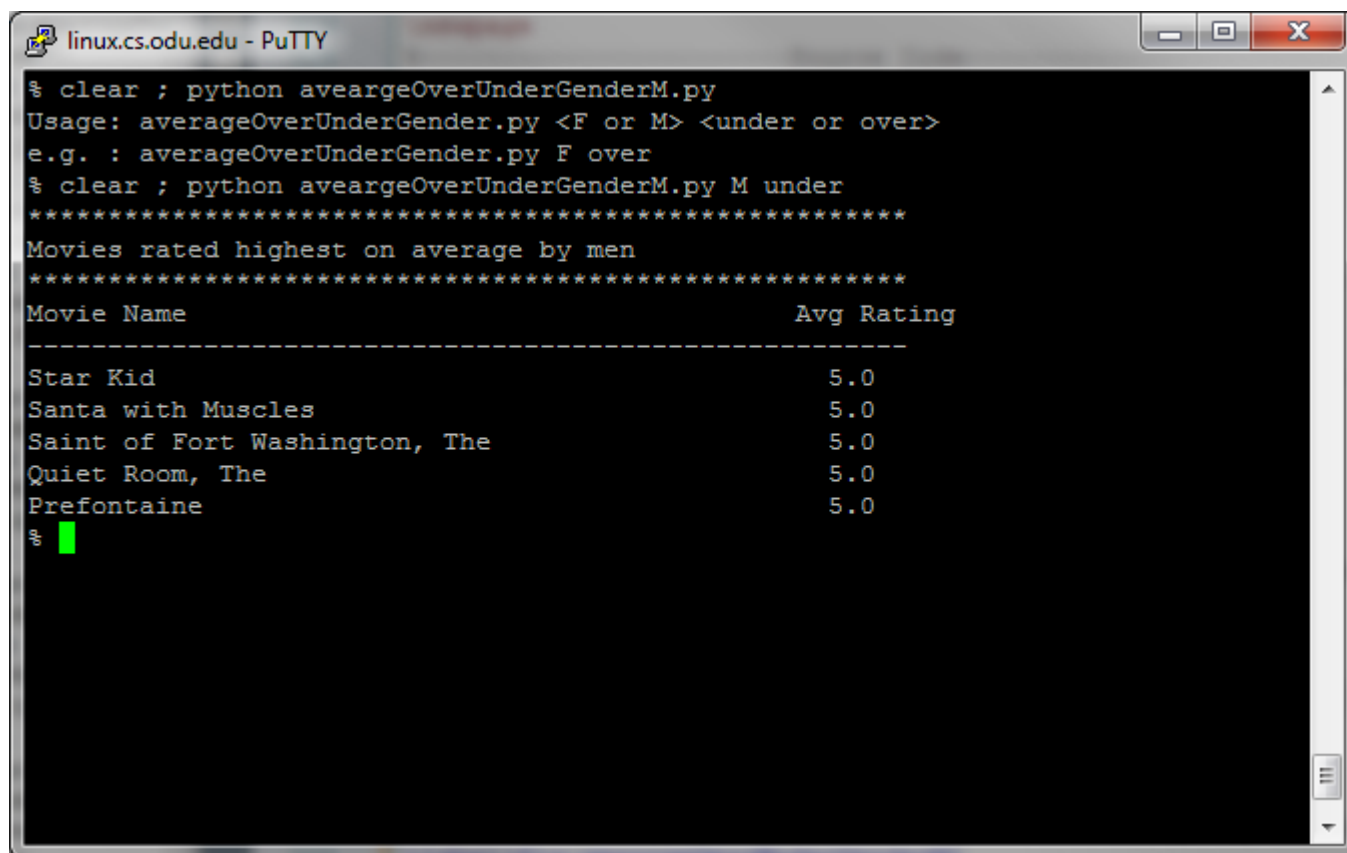
9.4.1 averageOverM.png



```
linux.cs.odu.edu - PuTTY
% clear ;python aveargeOverUnderGenderM.py
Usage: averageOverUnderGender.py <F or M> <under or over>
e.g. : averageOverUnderGender.py F over
%
%
% clear ; python aveargeOverUnderGenderM.py M over
*****
Movies rated highest on average by men
*****
Movie Name                               Avg Rating
-----
World of Apu, The (Apu Sansar)           5.0
Unstrung Heroes                          5.0
Two or Three Things I Know About Her     5.0
They Made Me a Criminal                  5.0
Strawberry and Chocolate (Fresa y chocolate) 5.0
% █
```

Figure 6: Movies rated highest on average by men over 40.

9.4.2 averageUnderM.png



```
linux.cs.odu.edu - PuTTY
% clear ; python aveargeOverUnderGenderM.py
Usage: averageOverUnderGender.py <F or M> <under or over>
e.g. : averageOverUnderGender.py F over
% clear ; python aveargeOverUnderGenderM.py M under
*****
Movies rated highest on average by men
*****
Movie Name                               Avg Rating
-----
Star Kid                                5.0
Santa with Muscles                      5.0
Saint of Fort Washington, The           5.0
Quiet Room, The                         5.0
Prefontaine                            5.0
% █
```

Figure 7: Movies rated highest on average by men under 40.

9.4.3 averageOverM.txt

```
1 % clear ; python averageOverUnderGender.py
2 Usage: averageOverUnderGender.py <F or M> <under or over>
3 e.g. : averageOverUnderGender.py F over
4 %
5 %
6 % clear ; python averageOverUnderGender.py M over
7 *****
8 Movies rated highest on average by men
9 *****
10 Movie Name                               Avg Rating
11 -----
12 World of Apu, The (Apu Sansar)           5.0
13 Unstrung Heroes                           5.0
14 Two or Three Things I Know About Her     5.0
15 They Made Me a Criminal                  5.0
16 Strawberry and Chocolate (Fresa y chocolate) 5.0
17 Star Kid                                 5.0
18 Spice World                              5.0
19 Solo                                     5.0
20 Rendezvous in Paris (Rendez-vous de Paris, Les) 5.0
21 Prefontaine                             5.0
22 Poison Ivy II                            5.0
23 Marlene Dietrich: Shadow and Light       5.0
24 Little Princess, The                    5.0
25 Little City                             5.0
26 Leading Man, The                        5.0
27 Late Bloomers                           5.0
28 Indian Summer                           5.0
29 Hearts and Minds                        5.0
30 Great Day in Harlem, A                  5.0
31 Grateful Dead                           5.0
32 Faithful                                5.0
33 Double Happiness                        5.0
34 Boxing Helena                           5.0
35 Aparajito                               5.0
36 Ace Ventura: When Nature Calls          5.0
37 Pather Panchali                          4.8
38 Whole Wide World, The                    4.7
39 A Chef in Love                           4.7
40 Close Shave, A                           4.6
41 Shanghai Triad (Yao a yao yao dao waipo qiao) 4.6
```

9.4.4 averageUnderM.txt

```
1 % clear ; python averageOverUnderGender.py
2 Usage: averageOverUnderGender.py <F or M> <under or over>
3 e.g. : averageOverUnderGender.py F over
4 %
5 %
6 % clear ; python averageOverUnderGender.py M under
7 *****
8 Movies rated highest on average by men
9 *****
10 Movie Name                               Avg Rating
11 -----
12 Star Kid                                5.0
13 Santa with Muscles                      5.0
14 Saint of Fort Washington, The           5.0
15 Quiet Room, The                        5.0
16 Prefontaine                            5.0
17 Perfect Candidate, A                   5.0
18 Maya Lin: A Strong Clear Vision         5.0
19 Magic Hour, The                        5.0
20 Love in the Afternoon                  5.0
21 Love Serenade                          5.0
22 Letter From Death Row, A              5.0
23 Leading Man, The                      5.0
24 Hugo Pool                             5.0
25 Entertaining Angels: The Dorothy Day Story 5.0
26 Delta of Venus                         5.0
27 Crossfire                             5.0
28 Angel Baby                            5.0
29 Aiqing wansui                          5.0
30 Winter Guest, The                     4.5
31 Two or Three Things I Know About Her   4.5
32 Sum of Us, The                        4.5
33 Sliding Doors                         4.5
34 Man of No Importance, A               4.5
35 Little Princess, The                 4.5
36 Innocents, The                       4.5
37 Grosse Fatigue                        4.5
38 Fille seule, La (A Single Girl)       4.5
39 Boy's Life 2                          4.5
40 Anna                                  4.5
41 Wallace & Gromit: The Best of Aardman Animation 4.5
```

10 Question 10:

What movie was rated highest on average by women over 40? By women under 40?

10.1 Approach

1. The approach is same as question 9.
2. To get the corresponding output for this question it will still use `averageOverUnderGender.py` but the only difference is changing the arguments.
3. Figure 8 is the output showing how to execute the program and displays the top 5 movies which are rated highest on average by men “over” 40.
4. Figure 9 is the output showing how to execute the program and displays the top 5 movies which are rated highest on average by men “under” 40.
5. The file `averageOverF.txt` and `averageUnderF.txt` displays movies which have the same average ratings.

10.2 Source Code

10.2.1 averageOverUnderGender.py

```
1 #!/usr/bin/env python
2
3
4 # program to get the highest average ratings by women under and over 40.
5 import sys
6 import operator
7 def main():
8     # Take the arguments from the command line
9     numOfArgs=len(sys.argv)
10    if numOfArgs<3 or numOfArgs>3:
11
12        print 'Usage: averageOverUnderGender.py <F or M> <under or over>'
13        print 'e.g. : averageOverUnderGender.py F over '
14        sys.exit(1)
15    gender          = sys.argv[1]
16    gender_under_over= sys.argv[2]
17
18    # Create a dictionary
19    movie_userid_avg = {}
20    movie_id_ratings = {}
21    movie_userid_avg = {}
22    users            = {}
23    count            = 1
24
25    # Read the files
26    readData = open('/home/bbokka/cs594/A8/u.data', 'r')
27    readItem = open('/home/bbokka/cs594/A8/u.item', 'r')
28    readUser = open('/home/bbokka/cs594/A8/u.user', 'r')
29
30    # Reading u.user file to get the male or female
31    for line in readUser:
32        split_input_line = line.strip().split('|')
33        user_id_from_user= split_input_line[0]
34        user_age         = int(split_input_line[1])
35        user_gender       = split_input_line[2]
36
37        if gender_under_over == 'over' :
38            if(gender.upper() == user_gender and user_age > 40):
39                users[user_id_from_user] = [user_age]
40
41        if gender_under_over == 'under' :
42            if(gender.upper() == user_gender and user_age < 40):
43                users[user_id_from_user] = [user_age]
44
45    readUser.close()
46
47    # reading the u.data file for item and ratings.
48    for line in readData:
49        split_input_line = line.strip().split('\t')
50        user_id_from_data = split_input_line[0]
51        movie_id_from_data= split_input_line[1]
52        movie_rating      = float(split_input_line[2])
53
```



```

54         if user_id_from_data in users:
55             try:
56                 movie_id_ratings[movie_id_from_data].append(movie_rating)
57             except KeyError:
58                 movie_id_ratings[movie_id_from_data] = list()
59                 movie_id_ratings[movie_id_from_data].append(movie_rating)
60     readData.close()
61
62     # Calculating the average.
63     for key in movie_id_ratings:
64         avg = sum(movie_id_ratings[key])/len(movie_id_ratings[key])
65         movie_userid_avg[key] = [float(avg)]
66
67     # Reading the u.item file for movie name.
68     for each_line in readItem:
69         split_each_line = each_line.strip().split('|')
70
71         movie_item_id = split_each_line[0]
72         movie_name_split = split_each_line[1].split('(1')
73         movie_name = movie_name_split[0]
74
75         try:
76             movie_userid_avg[movie_item_id].append(movie_name)
77         except KeyError:
78             pass
79     readItem.close()
80     print '*' * 55
81     print "Movies rated highest on average by women"
82     print "*" * 55
83     print "Movie Name\t\t\t\t\t","Avg Rating"
84     print "-" * 55
85     # sorting the movies from highest to lowest based on the average value.
86     for key, value in sorted(movie_userid_avg.items(), key=lambda e: e[1], reverse=
87     True):
88         if (count <=5 ):
89             print '{:<50}{:<0.1f} '.format(value[1], value[0])
90             count += 1
91
92 if __name__ == "__main__":
93     try:
94         main()
95     except KeyboardInterrupt:
96         sys.exit(1)

```

10.3 Input Files

10.3.1 u.data

	userid	itemid	rating	timestamp
1				
2				
3	196	242	3	881250949
4	186	302	3	891717742
5	22	377	1	878887116
6	244	51	2	880606923
7	166	346	1	886397596
8	298	474	4	884182806
9	115	265	2	881171488

10.3.2 u.item

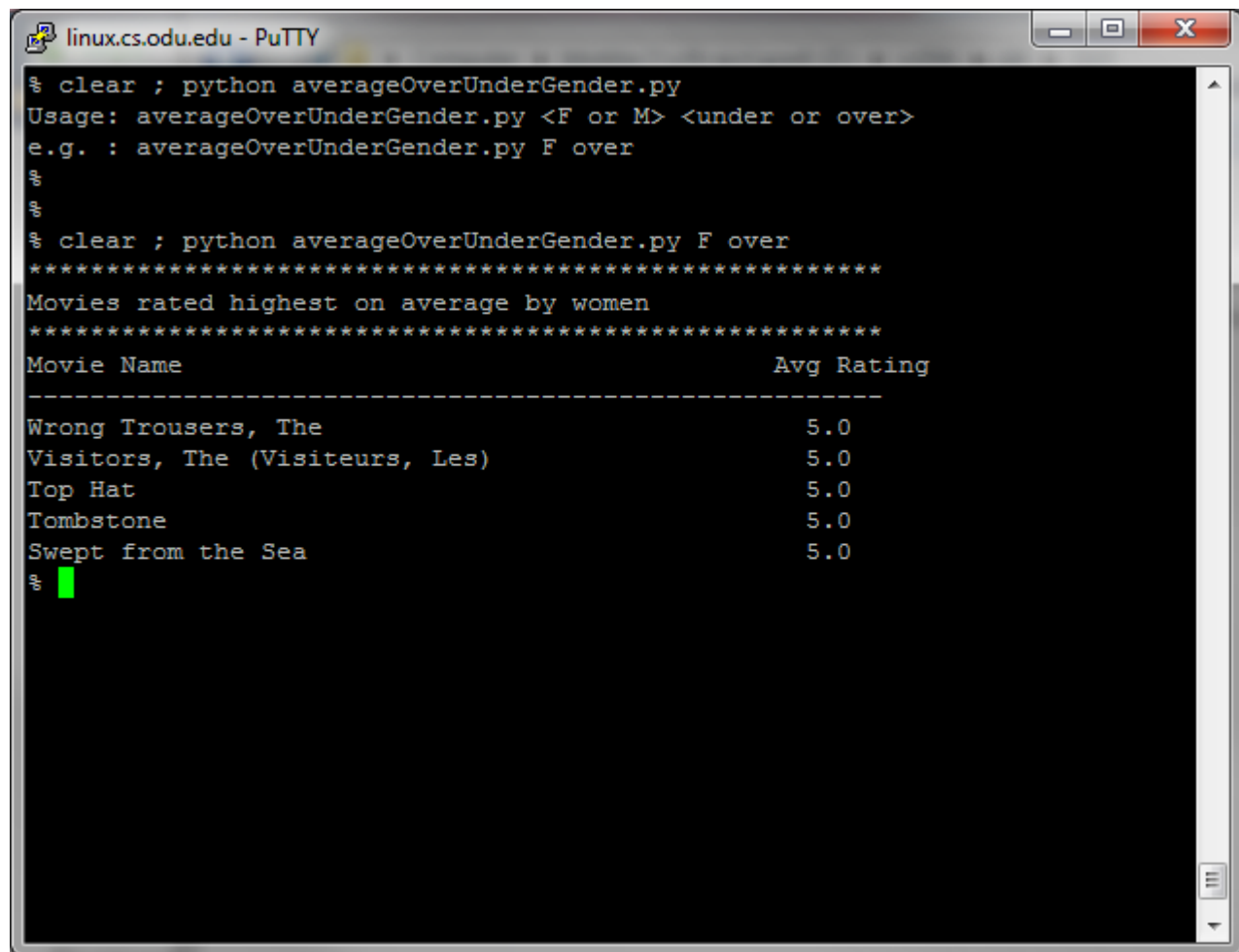
[illegible]

10.3.3 u.user

	userid	age	gender	occupation	zip_code
1	1	24	M	technician	85711
2	2	53	F	other	94043
3	3	23	M	writer	32067
4	4	24	M	technician	43537
5	5	33	F	other	15213
6	6	42	M	executive	98101
7	7	57	M	administrator	91344

10.4 Output Files

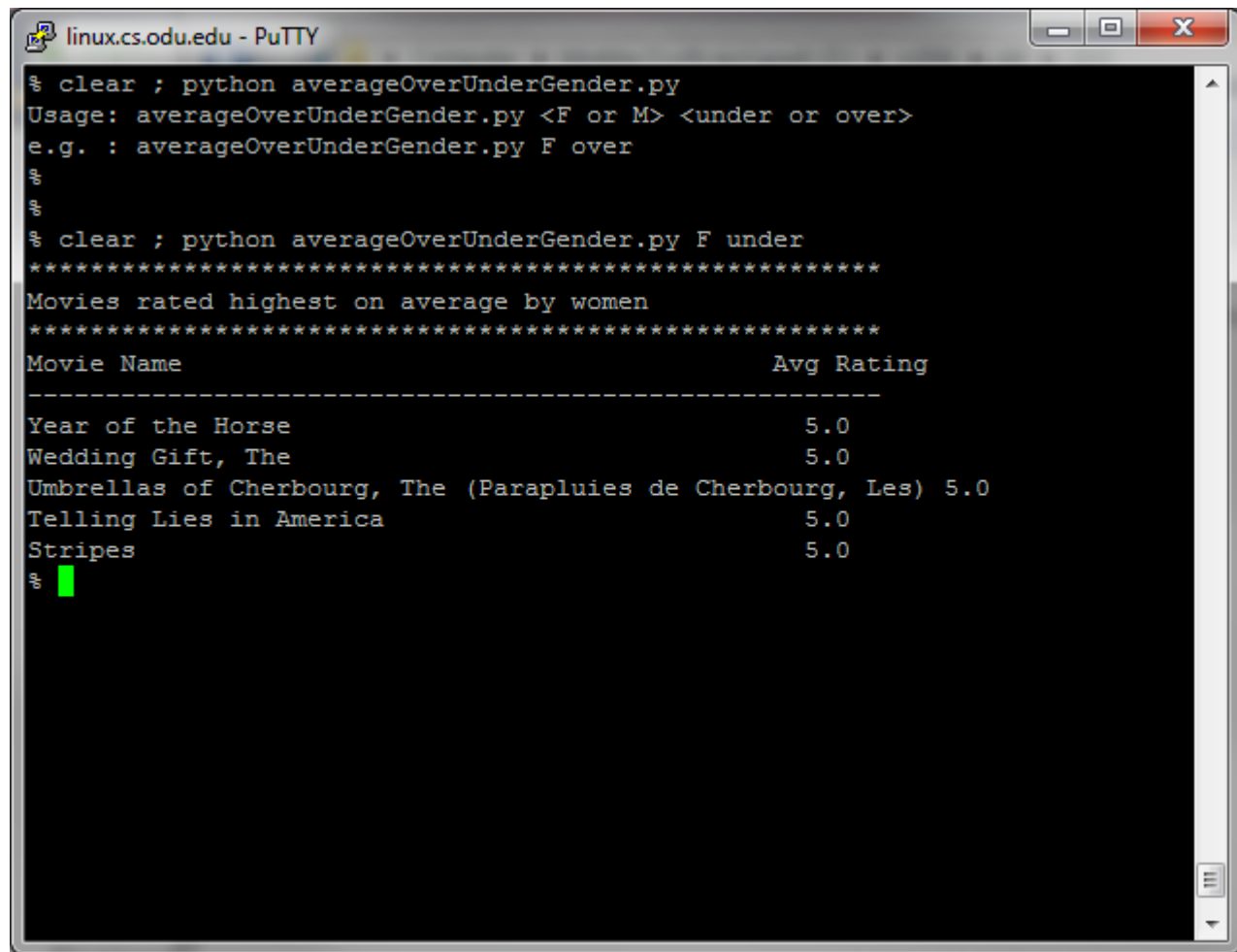
10.4.1 averageOverF.png



```
linux.cs.odu.edu - PuTTY
% clear ; python averageOverUnderGender.py
Usage: averageOverUnderGender.py <F or M> <under or over>
e.g. : averageOverUnderGender.py F over
%
%
% clear ; python averageOverUnderGender.py F over
*****
Movies rated highest on average by women
*****
Movie Name                                Avg Rating
-----
Wrong Trousers, The                      5.0
Visitors, The (Visiteurs, Les)          5.0
Top Hat                                 5.0
Tombstone                               5.0
Swept from the Sea                      5.0
% █
```

Figure 8: Movies rated highest on average by women over 40.

10.4.2 averageUnderF.png



```
linux.cs.odu.edu - PuTTY
% clear ; python averageOverUnderGender.py
Usage: averageOverUnderGender.py <F or M> <under or over>
e.g. : averageOverUnderGender.py F over
%
%
% clear ; python averageOverUnderGender.py F under
*****
Movies rated highest on average by women
*****
Movie Name                               Avg Rating
-----
Year of the Horse                         5.0
Wedding Gift, The                         5.0
Umbrellas of Cherbourg, The (Parapluies de Cherbourg, Les) 5.0
Telling Lies in America                   5.0
Stripes                                  5.0
% █
```

Figure 9: Movies rated highest on average by women under 40.

10.4.3 averageOverF.txt

```
1 % clear ; python averageOverUnderGender.py
2 Usage: averageOverUnderGender.py <F or M> <under or over>
3 e.g. : averageOverUnderGender.py F over
4 %
5 %
6 % clear ; python averageOverUnderGender.py F over
7 *****
8 Movies rated highest on average by women
9 *****
10 Movie Name                               Avg Rating
11 -----
12 Wrong Trousers, The                      5.0
13 Visitors, The (Visiteurs, Les)          5.0
14 Top Hat                                  5.0
15 Tombstone                                5.0
16 Swept from the Sea                      5.0
17 Shallow Grave                           5.0
18 Shall We Dance?                         5.0
19 Safe                                     5.0
20 Quest, The                              5.0
21 Pocahontas                              5.0
22 Nightmare Before Christmas, The         5.0
23 Mina Tannenbaum                          5.0
24 Mary Shelley's Frankenstein             5.0
25 Ma vie en rose (My Life in Pink)        5.0
26 Letter From Death Row, A               5.0
27 In the Bleak Midwinter                  5.0
28 Great Dictator, The                     5.0
29 Grand Day Out, A                        5.0
30 Gold Diggers: The Secret of Bear Mountain 5.0
31 Funny Face                              5.0
32 Foreign Correspondent                    5.0
33 Bride of Frankenstein                   5.0
34 Best Men                                5.0
35 Band Wagon, The                         5.0
36 Balto                                   5.0
37 Angel Baby                              5.0
38 Once Were Warriors                      4.8
39 Fantasia                                4.7
40 Last of the Mohicans, The              4.7
41 Christmas Carol, A                     4.6
```

10.4.4 averageUnderF.txt

```
1 % clear ; python averageOverUnderGender.py
2 Usage: averageOverUnderGender.py <F or M> <under or over>
3 e.g. : averageOverUnderGender.py F over
4 %
5 %
6 % clear ; python averageOverUnderGender.py F under
7 *****
8 Movies rated highest on average by women
9 *****
10 Movie Name                               Avg Rating
11 -----
12 Year of the Horse                         5.0
13 Wedding Gift , The                      5.0
14 Umbrellas of Cherbourg, The (Parapluies de Cherbourg, Les) 5.0
15 Telling Lies in America                 5.0
16 Stripes                                 5.0
17 Someone Else 's America                 5.0
18 Prefontaine                            5.0
19 Nico Icon                               5.0
20 Mina Tannenbaum                         5.0
21 Maya Lin: A Strong Clear Vision         5.0
22 Horseman on the Roof, The (Hussard sur le toit , Le) 5.0
23 Heaven 's Prisoners                     5.0
24 Grace of My Heart                       5.0
25 Faster Pussycat! Kill! Kill!           5.0
26 Everest                                5.0
27 Don't Be a Menace to South Central While Drinking Your Juice in the Hood 5.0
28 Backbeat                               5.0
29 Wallace & Gromit: The Best of Aardman Animation 4.8
30 Paradise Lost: The Child Murders at Robin Hood Hills 4.8
31 Anne Frank Remembered                  4.8
32 Shawshank Redemption, The              4.7
33 Shall We Dance?                        4.7
34 Close Shave, A                         4.7
35 Schindler 's List                       4.7
36 Margaret 's Museum                     4.7
37 Women, The                             4.5
38 To Catch a Thief                       4.5
39 Thin Man, The                           4.5
40 Some Folks Call It a Sling Blade       4.5
41 Primary Colors                         4.5
```

References

- [1] Data files used for assignment. <http://files.grouplens.org/datasets/movielens/ml-100k/>, April 1998. Data collected by GroupLens Research Project.
- [2] K. Arthur Endsley. recommendations.py program. <https://github.com/arthur-e/Programming-Collective-Intelligence/blob/master/chapter2/recommendations.py>, December 2012.
- [3] Frederick Hamidi. Python: Run function from the command line. <http://stackoverflow.com/questions/3987041/python-run-function-from-the-command-line>, October 2010.
- [4] Devin Jeanpierre. Sort a python dictionary by value. <http://stackoverflow.com/questions/613183/sort-a-python-dictionary-by-value>, March 2009.
- [5] Toby Segaran. *Programming Collective Intelligence*. O'Reilly Media, August 2007.

□