## HW#7 - Recommendation Systems
PRASHANT TOMAR
04/09/2023

# Q1 - Find 3 and Movie Preference

## Answer

Initially, we will identify three users from the u.user.txt dataset that have similar characteristics to my age, gender, and occupation. The identified users will be saved in a distinct list. To convert the data into a DataFrame, we will use the same code from Chapter 2 of Programming Collective Intelligence.

Subsequently, we will loop through each user and obtain their ratings data from the u.data.txt file. For this operation and the rest of the task, we will use pandas DataFrame to filter and index data.

During each iteration, we will extract movie ratings provided by the user and save them in a different DataFrame, which will then be sorted according to the rating. We will then use the iloc function of the DataFrame to select and save the top three and last three rows in separate lists.

```python
import pandas as pd

def load_item_dataset():
    # Define column names for u.item file
    col = ['movie_id','movie_title','release_date','video_release_date'
    ,'IMDb_URL','unknown','Action','Adventure','Animation','Childrens','
    Comedy','Crime','Documentary','Drama','Fantasy','Film_Noir','Horror'
    ,'Musical','Mystery','Romance','Sci_Fi','Thriller','War','Western']

    # Read u.item file into item_dataset DataFrame
    item_dataset = pd.read_csv('u.item.txt', delimiter='|', header=None
    , names=col)

    return item_dataset


def find_similar_users(age, gender, occupation):
    similar_user = []
    for line in open('u.user.txt'):
        (user_id, user_age, user_gender, user_occ, user_zip) = line.
    split('|')
        if user_age == age and user_gender == gender and user_occ ==
```

```python
       occupation:
18             similar_user.append(user_id)
19
20     return similar_user
21
22
23 def load_rating_dataset(item_dataset):
24     # Create a list of dictionaries containing rating data for each
   item
25     rating_item_list = []
26     for line in open('u.data.txt'):
27         (user_id, item_id, rating, timestamp) = line.split('\t')
28
29         # Filter item_dataset to only include the current item
30         item = item_dataset[item_dataset['movie_id'] == pd.to_numeric(
   item_id)]
31
32         # Add current rating data to rating_item_list as a dictionary
33         rating_item = {'user_id': user_id, 'item_id': item_id, '
   item_name': item.iloc[0,1], 'rating': rating, 'timestamp': timestamp
   }
34         rating_item_list.append(rating_item)
35
36     # Convert rating_item_list to a DataFrame
37     rating_dataset = pd.DataFrame(rating_item_list)
38
39     return rating_dataset
40
41
42 def find_favorite_and_least_fav_films(similar_user, rating_dataset):
43     favorite_films_list = []
44     least_fav_films_list = []
45
46     for user in similar_user[0:3]:
47         user_rated_data = rating_dataset[rating_dataset['user_id'] ==
   user]
48         user_rated_data = user_rated_data.sort_values('rating',
   ascending=False)
49         user_rated_data_fav_films = user_rated_data.iloc[0:3, :]
50         user_rated_data_least_fav_films = user_rated_data.iloc[-3:]
51         favorite_films_list.append(user_rated_data_fav_films)
52         least_fav_films_list.append(user_rated_data_least_fav_films)
53
54     return favorite_films_list, least_fav_films_list
55
56
57 if __name__ == '__main__':
```

```
58    # Load item dataset
59    item_dataset = load_item_dataset()
60
61    # Find similar users
62    similar_user = find_similar_users('30', 'M', 'student')
63
64    # Load rating dataset
65    rating_dataset = load_rating_dataset(item_dataset)
66
67    # Find favorite and least favorite films for each user in
      similar_user
68    favorite_films_list, least_fav_films_list =
      find_favorite_and_least_fav_films(similar_user, rating_dataset)
69
70    # Print favorite and least favorite films for each user in
      similar_user
71    for ds in favorite_films_list:
72        print(ds.to_string(index=False))
73
74    for lds in least_fav_films_list:
75        print(lds.to_string(index=False))
```

**Listing 1:** Finding users with my preferences.

Below are the result of top favorite and least favorite movies of the three similar preferences users.

**Table 1:** 774 User Favorite Movies

| User | Movie id | Movie Title | Rating |
|------|----------|-------------|--------|
| 774 | 180 | Apocalypse Now (1979) | 5 |
| 774 | 202 | Groundhog Day (1993) | 5 |
| 774 | 193 | Right Stuff, The (1983) | 5 |

**Table 2:** 869 User Favorite Movies

| User | Movie id | Movie Title | Rating |
|------|----------|-------------|--------|
| 869 | 412 | Very Brady Sequel, A (1996) | 5 |
| 869 | 151 | Willy Wonka and the Chocolate Factory (1971) | 5 |
| 869 | 127 | Godfather, The (1972) | 5 |

Table 3: 774 User Least Favorite Movies

| User | Movie id | Movie Title | Rating |
|------|----------|-------------|--------|
| 774 | 411 | Nutty Professor, The (1996) | 1 |
| 774 | 168 | Monty Python and the Holy Grail (1974) | 1 |
| 774 | 393 | Mrs. Doubtfire (1993) | 1 |

Table 4: 869 User Least Favorite Movies

| User | Movie id | Movie Title | Rating |
|------|----------|-------------|--------|
| 869 | 1061 | Evening Star, The (1996) | 1 |
| 869 | 476 | First Wives Club, The (1996) | 1 |
| 869 | 284 | Tin Cup (1996) | 1 |

## Analysis

Out of all the users, I will choose user 774 as my "substitute" because I have watched Apocalypse Now and it happens to be one of my favorite movies. After searching for the other movies that were present in different users' favorite lists, I found out that they were mainly non-realistic and fictional movies, which I am not particularly fond of.

# Q2 - Most Correlated User  Least Correlated User with respect to Substitute You.

## Answer

To accomplish this, we will utilize the code provided in Chapter 2 of Programming Collective Intelligence. The code consists of two key functions, one that computes the distance between two users and another that identifies the top match based on this distance function. We will use both Euclidean Distance and Pearson method to calculate the distance between users and find the most correlated user. However, we will only share the output generated using Pearson method since it is the most reliable and effective method.

```python
import numpy as np
import math

def sim_distance(prefs: dict[str, dict[str, float]], person1: str,
    person2: str) -> float:
    '''
     Returns a distance-based similarity score for person1 and person2.
    '''
    shared_items = {item for item in prefs[person1] if item in prefs[
    person2]}
    if len(shared_items) == 0:
        return 0
    sum_of_squares = sum((prefs[person1][item] - prefs[person2][item])
    ** 2 for item in shared_items)
    return 1 / (1 + math.sqrt(sum_of_squares))

def sim_pearson(prefs: dict[str, dict[str, float]], person1: str,
    person2: str) -> float:
    '''
     Returns the Pearson correlation coefficient for person1 and person2
    .
    '''
    shared_items = {item for item in prefs[person1] if item in prefs[
    person2]}
    if len(shared_items) == 0:
        return 0
    n = len(shared_items)
    sum1 = sum(prefs[person1][item] for item in shared_items)
    sum2 = sum(prefs[person2][item] for item in shared_items)
    sum1_sq = sum(prefs[person1][item] ** 2 for item in shared_items)
    sum2_sq = sum(prefs[person2][item] ** 2 for item in shared_items)
    p_sum = sum(prefs[person1][item] * prefs[person2][item] for item in
    shared_items)
```

```python
27      num = p_sum - (sum1 * sum2) / n
28      den = math.sqrt((sum1_sq - (sum1 ** 2) / n) * (sum2_sq - (sum2 **
    2) / n))
29      if den == 0:
30          return 0
31      r = num / den
32      return r
33
34  def top_matches(
35      prefs: dict[str, dict[str, float]],
36      person: str,
37      n: int = 5,
38      similarity: callable = sim_pearson,
39  ) -> list[tuple[float, str]]:
40      '''
41      Returns the top matches for person from the prefs dictionary.
42      Number of results and similarity function are optional params.
43      '''
44      scores = [(similarity(prefs, person, other), other) for other in
    prefs if other != person]
45      scores.sort(reverse=True)
46      return scores[:n]
47
48  def least_matches(
49      prefs: dict[str, dict[str, float]],
50      person: str,
51      n: int = 5,
52      similarity: callable = sim_pearson,
53  ) -> list[tuple[float, str]]:
54      '''
55      Returns the least matches for person from the prefs dictionary.
56      Number of results and similarity function are optional params.
57      '''
58      scores = [(similarity(prefs, person, other), other) for other in
    prefs if other != person]
59      scores.sort()
60      return scores[:n]
61
62
63  def load_movie_lens():
64      # Get movie titles
65      with open('u.item.txt', encoding='utf-8') as f:
66          movies = {}
67          for line in f:
68              id, title = line.split('|')[:2]
69              movies[id] = title
70
```

```python
71    # Load data
72    with open('u.data.txt', encoding='utf-8') as f:
73        prefs = {}
74        for line in f:
75            user, movie_id, rating, _ = line.split('\t')
76            prefs.setdefault(user, {})
77            prefs[user][movies[movie_id]] = float(rating)
78
79    return prefs
80
81
82 prefs = load_movie_lens()
83
84 list_user = top_matches(prefs, '774', n=5)
85 print(list_user)
86
87 list_user = top_matches(prefs, '774', n=5, similarity=sim_distance)
88 print(list_user)
89
90 list_user = least_matches(prefs, '774', n=5)
91 print(list_user)
92
93 list_user = least_matches(prefs, '774', n=5, similarity=sim_distance)
94 print(list_user)
```

**Listing 2:** Finding users with my preferences.

It can be observed that I have modified the topMatch function into a LeastMatch function by altering the sorting code to maintain the order instead of reversing it. I will use both the topMatch and leastMatch functions, with user 774 as my substitute, to identify the users that are most and least correlated to my substitute.

The results for the top match and least match users of my substitute are provided below.

**Table 5:** Top Match Correlated User

| User | Correlation Result |
|------|--------------------|
| 428 | 1.0000000000000033 |
| 300 | 1.0000000000000007 |
| 170 | 1.0000000000000007 |
| 909 | 1.0 |
| 762 | 1.0 |

**Table 6:** Top Match Correlated User

| User | Correlation Result |
|------|--------------------|
| 284  | -1.0000000000000007 |
| 107  | -1.0 |
| 155  | -1.0 |
| 240  | -1.0 |
| 252  | -1.0 |

## Analysis

The following users have been identified as being more similar to my substitute, and their movie preferences have been analyzed. Although some of the movie genres aligned with my interests, others did not. I attempted to use Euclidean Distance to find correlated users as well, and while some users matched with Pearson, not all of them did.

# Q3 - Compute Ratings, Top 5 and Bottom 5 recommendations for movies

## Answer

We will utilize the code provided in Chapter 2 of the book Programming Collective Intelligence. Specifically, we will employ the getRecommendation function, which utilizes the Pearson function to identify the movies that my substitute should watch based on the closest distance.

The implementation of the current process is presented below.

```python
from numpy import sqrt

def distance_similarity(prefs, p1, p2):
    si = {}
    for item in prefs[p1]:
        if item in prefs[p2]:
            si[item] = 1
    if len(si) == 0:
        return 0
    sum_of_squares = sum([pow(prefs[p1][item] - prefs[p2][item], 2) for item in
                          prefs[p1] if item in prefs[p2]])
    return 1 / (1 + sqrt(sum_of_squares))

def pearson_similarity(prefs, p1, p2):
    si = {}
    for item in prefs[p1]:
        if item in prefs[p2]:
            si[item] = 1
    if len(si) == 0:
        return 0
    n = len(si)
    sum1 = sum([prefs[p1][it] for it in si])
    sum2 = sum([prefs[p2][it] for it in si])
    sum1Sq = sum([pow(prefs[p1][it], 2) for it in si])
    sum2Sq = sum([pow(prefs[p2][it], 2) for it in si])
    pSum = sum([prefs[p1][it] * prefs[p2][it] for it in si])
    num = pSum - sum1 * sum2 / n
    den = sqrt((sum1Sq - pow(sum1, 2) / n) * (sum2Sq - pow(sum2, 2) / n))
    if den == 0:
        return 0
    r = num / den
    return r
```

```python
34 def get_recommendations(prefs, person, similarity=pearson_similarity):
35     totals = {}
36     sim_sums = {}
37     for other in prefs:
38         if other == person:
39             continue
40         sim = similarity(prefs, person, other)
41         if sim <= 0:
42             continue
43         for item in prefs[other]:
44             if item not in prefs[person] or prefs[person][item] == 0:
45                 totals.setdefault(item, 0)
46                 totals[item] += prefs[other][item] * sim
47                 sim_sums.setdefault(item, 0)
48                 sim_sums[item] += sim
49     rankings = [(total / sim_sums[item], item) for (item, total) in
50                 totals.items()]
51     rankings.sort()
52     rankings.reverse()
53     return rankings
54
55 def load_movie_lens():
56     movies = {}
57     for line in open('u.item.txt'):
58         (id, title) = line.split('|')[0:2]
59         movies[id] = title
60     prefs = {}
61     for line in open('u.data.txt'):
62         (user, movieid, rating, ts) = line.split('\t')
63         prefs.setdefault(user, {})
64         prefs[user][movies[movieid]] = float(rating)
65     return prefs
66
67 prefs = load_movie_lens()
68 list_user = get_recommendations(prefs, '774')
69
70 print(list_user[0:5])
71
72 print(list_user[-5:])
```

**Listing 3:** Recommendation movies code

After receiving the list of suggested movies, we will store the recommendations in a list and select the top 5 and bottom 5 movies. This will allow us to display the movies that are highly recommended as well as those that are least recommended. The results of this process are presented below.

**Table 7:** Top 5 recommended movies

| Movies | Rating |
| --- | --- |
| Prefontaine (1997) | 5.0 |
| Tough and Deadly (1995) | 5.0 |
| They Made Me a Criminal (1939) | 5.0 |
| Star Kid (1997) | 5.0 |
| Someone Else's America (1995) | 5.0 |

**Table 8:** Least 5 recommended movies

| Movies | Rating |
| --- | --- |
| Amityville: A New Generation (1993) | 1.0 |
| Amityville Curse, The (1990) | 1.0 |
| Amityville 3-D (1983) | 1.0 |
| Amityville 1992: It's About Time | 1.0 |
| 3 Ninjas: High Noon At Mega Mountain | 1.0 |

## Analysis

Upon reviewing the top 5 recommended movies, I searched for some of them and found that most of them seem intriguing, and I assume they will be good. However, after examining the least recommended movies, I found that none of them appear interesting to me at all.

# Q4 - Correlated Movies.

## Answer

We will utilize the code we used in Q1 and add a new function that will transform the movie data into a format similar to the user data. This will enable us to compute the distance between movies and identify the movies that are most correlated.

The implementation of this process is presented below.

```python
1  import numpy as np
2
3  def sim_distance(prefs, p1, p2):
4      si = {}
5      for item in prefs[p1]:
6          if item in prefs[p2]:
7              si[item] = 1
8      if not si:
9          return 0
10     sum_of_squares = sum([pow(prefs[p1][item] - prefs[p2][item], 2) for
       item in si])
11     return 1 / (1 + np.sqrt(sum_of_squares))
12
13 def sim_pearson(prefs, p1, p2):
14     si = {}
15     for item in prefs[p1]:
16         if item in prefs[p2]:
17             si[item] = 1
18     if not si:
19         return 0
20     n = len(si)
21     sum1 = sum([prefs[p1][it] for it in si])
22     sum2 = sum([prefs[p2][it] for it in si])
23     sum1_sq = sum([pow(prefs[p1][it], 2) for it in si])
24     sum2_sq = sum([pow(prefs[p2][it], 2) for it in si])
25     p_sum = sum([prefs[p1][it] * prefs[p2][it] for it in si])
26     num = p_sum - (sum1 * sum2 / n)
27     den = np.sqrt((sum1_sq - pow(sum1, 2) / n) * (sum2_sq - pow(sum2,
       2) / n))
28     if den == 0:
29         return 0
30     r = num / den
31     return r
32
33 def top_matches(prefs, person, n=5, similarity=sim_pearson):
34     scores = [(similarity(prefs, person, other), other) for other in
```

```
        prefs if other != person]
35      scores.sort(reverse=True)
36      return scores[:n]
37
38 def least_matches(prefs, person, n=5, similarity=sim_pearson):
39      scores = [(similarity(prefs, person, other), other) for other in
        prefs if other != person]
40      scores.sort()
41      return scores[:n]
42
43 def transform_prefs(prefs):
44      result = {}
45      for person in prefs:
46          for item in prefs[person]:
47              result.setdefault(item, {})
48              result[item][person] = prefs[person][item]
49      return result
50
51 def load_movie_lens():
52      movies = {}
53      with open('u.item.txt', encoding='ISO-8859-1') as f:
54          for line in f:
55              (id, title) = line.split('|')[0:2]
56              movies[id] = title
57      prefs = {}
58      with open('u.data.txt', encoding='ISO-8859-1') as f:
59          for line in f:
60              (user, movieid, rating, ts) = line.split('\t')
61              prefs.setdefault(user, {})
62              prefs[user][movies[movieid]] = float(rating)
63      return prefs
64
65 prefs = load_movie_lens()
66 mod_prefs = transform_prefs(prefs)
67
68 list_movies = top_matches(mod_prefs, 'Silence of the Lambs, The (1991)'
    , n=5)
69 print(list_movies)
70
71 list_movies = least_matches(mod_prefs, 'Silence of the Lambs, The
    (1991)', n=5)
72 print(list_movies)
73
74 list_movies = top_matches(mod_prefs, 'Pulp Fiction (1994)', n=5)
75 print(list_movies)
76
77 list_movies = least_matches(mod_prefs, 'Pulp Fiction (1994)', n=5)
```

```
78 print(list_movies)
```

**Listing 4:** Further tweet analysis

As you can see, I have created a new function called leastMatch by modifying the topMatch function and removing the reverse functionality. This function allows us to identify the least correlated movies and most correlated movies. Once we have obtained the list of movies from these functions, we will display them. Additionally, we are sending n=5 into the function to specify the number of data that we require.

The result of this process is presented below.

**Table 9:** Most correlated movies of my Favorite Movies

| Movies | Correlation Values |
| --- | --- |
| Mrs. Dalloway (1997) | 1.0 |
| One Night Stand (1997) | 1.0 |
| Smile Like Yours, A (1997) | 1.0 |
| Year of the Horse (1997) | 1.0 |
| Total Eclipse (1995) | 1.0 |

**Table 10:** Least correlated movies of my Favorite Movies

| Movies | Correlation Values |
| --- | --- |
| Children of the Revolution (1996) | -1.0 |
| My Favorite Season (1993) | -1.0 |
| Half Baked (1998) | -1.0 |
| Broken English (1996) | -1.0 |
| Caro Diario (Dear Diary) (1994) | -1.0 |

**Table 11:** Most correlated movies of my Least Favorite Movies

| Movies | Correlation Values |
| --- | --- |
| 'Maya Lin: A Strong Clear Vision (1994) | 1.0 |
| 'Crossfire (1947) | 1.0 |
| 'Zeus and Roxanne (1997) | 1.0 |
| 'Wedding Gift, The (1994) | 1.0 |
| 'Two Deaths (1995) | 1.0 |

**Table 12:** Least correlated movies of my Least Favorite Movies

| Movies | Correlation Values |
|---|---|
| 'Mr. Magoo (1997) | -1.0 |
| 'Telling Lies in America (1997) | -1.0 |
| 'American Dream (1990) | -1.0 |
| Anna (1996) | -1.0 |
| 'Bewegte Mann, Der (1994) | -1.0 |

## Analysis

I am satisfied with the results as the recommended movies are of my liking and interest. I researched the movies on the internet and read their plot to confirm my preferences.

I am not a fan of movies that contain unrealistic elements, such as Godzilla, and this can be seen in the least favorite most correlated table where there is a Mr. Magoo movie that is related to fictional stories that I am not interested in. This demonstrates that the Pearson method of calculating the distance is effective in determining the correlation between objects.

# Q6 - Rank all 1682 movies

## Answer

The process for this task involves loading both the movie data and movie rating data, similarly to the previous example. After loading the data, each movie will be looped through, and the total count of ratings and average rating will be determined.

Once all the necessary data is obtained, a separate DataFrame will be created to perform further calculations on. A new column will be added to this DataFrame, which calculates the rating weight by multiplying the count of votes with the average rating.

Implementation of the current process is shared below,

```python
1  import pandas as pd
2  from pandas import DataFrame
3  import numpy as np
4
5  col = ['movie_id','movie_title','release_date','video_release_date','
       IMDb_URL','unknown','Action','Adventure','Animation','Childrens','
       Comedy','Crime','Documentary','Drama','Fantasy','Film_Noir','Horror'
       ,'Musical','Mystery','Romance','Sci_Fi','Thriller','War','Western']
6  item_dataset = pd.read_csv('u.item.txt', delimiter = "|", header=None,
       names=col)
7
8  rating_item_list = []
9  for line in open('u.data.txt'):
10     (user_id, item_id, rating, timestamp) = line.split('\t')
11
12     rating_item = {'user_id': user_id, 'item_id': item_id,
13                     'rating': rating, 'timestamp': timestamp
14                     }
15     rating_item_list.append(rating_item)
16
17 rating_dataset = DataFrame(rating_item_list)
18 rating_dataset['rating'] = pd.to_numeric(rating_dataset['rating'])
19
20 rating_list = []
21 for movie in item_dataset.iterrows():
22     ratings = rating_dataset[rating_dataset['item_id'] == str(movie[1][
       'movie_id'])]
23     average_rating = round(np.sum(ratings.iloc[:, 2]) / len(ratings.
       index), 2)
24     rating_count = len(ratings.index)
25
26     rating_object = {
```

```
27          'movie_id' : movie[1]['movie_id'],
28          'movie_title': movie[1]['movie_title'],
29          'average_rating': average_rating,
30          'rating_count': rating_count,
31          'weighted_rating' : average_rating * rating_count
32      }
33    rating_list.append(rating_object)
34
35 rating_dataset = DataFrame(rating_list)
36
37 rating_dataset = rating_dataset.sort_values('weighted_rating',
       ascending=False)
38 rating_dataset.to_csv('rating_sorted_dataset.csv')
39 print(rating_dataset.iloc[0:10,:])
40 print(rating_dataset.iloc[-10:])
```

**Listing 5:** Ranking the movies

After loading the movie and movie rating data, we will iterate through each movie to calculate the total count of ratings and the average rating. Then, we will create a separate DataFrame and calculate the rating weight by multiplying the count of votes with the average rating.

After calculating the rating weight, we will sort the DataFrame based on the weighted average rating column and retrieve the top and bottom 10 using the iloc function of DataFrame. Finally, we will save the ranked results in a separate file.

The results are shown below.

**Table 13:** Top 10 ranked movies

| Movies | Avg. Rating | Count | Weighted Average |
|---|---|---|---|
| Star Wars (1977)) | 4.36 | 583 | 2541.88 |
| Fargo (1996) | 4.16 | 508 | 2113.28 |
| Return of the Jedi (1983) | 4.01 | 507 | 2033.07 |
| Contact (1997) | 3.80 | 509 | 1934.20 |
| Raiders of the Lost Ark (1981) | 4.25 | 420 | 1785.00 |
| Godfather, The (1972) | 4.28 | 413 | 1767.64 |
| English Patient, The (1996) | 3.66 | 481 | 1760.46 |
| Toy Story (1995) | 3.88 | 452 | 1753.76 |
| Silence of the Lambs, The (1991) | 4.29 | 390 | 1673.10 |
| Scream (1996) | 3.44 | 478 | 1644.32 |

**Table 14:** Least 10 ranked movies

| Movies | Avg. Rating | Count | Weighted Average |
|---|---|---|---|
| Hungarian Fairy Tale, A (1987) | 1.0 | 1 | 1.0 |
| Pharaoh's Army (1995) | 1.0 | 1 | 1.0 |
| Modern Affair, A (1995) | 1.0 | 1 | 1.0 |
| Wend Kuuni (God's Gift) (1982) | 1.0 | 1 | 1.0 |
| Touki Bouki (Journey of the Hyena) (1973) | 1.0 | 1 | 1.0 |
| Woman in Question, The (1950) | 1.0 | 1 | 1.0 |
| Quartier Mozart (1992) | 1.0 | 1 | 1.0 |
| Girl in the Cadillac (1995) | 1.0 | 1 | 1.0 |
| Nobody Loves Me (Keiner liebt mich) (1994) | 1.0 | 1 | 1.0 |
| Liebelei (1933) | 1.0 | 1 | 1.0 |

## Analysis

The purpose of multiplying the count with the average rating is to calculate the weighted average rating of the movies, which makes it easier to rank them. For instance, in table 13, even though The Godfather has a higher average rating of 4.28, it has only 413 user ratings compared to Fargo, which has a lower average rating of 4.16 but a higher count of user ratings at 508. That's why Fargo is ranked higher than The Godfather.

# Q7 - Rank the same 1682 movies according to today's IMDB data

## Answer

To complete this task, we will start by downloading two files, basic.csv and rating.csv, from today's Imdb movies data. Then, we will use pandas' read_csv function to load these CSV files. After loading the files, we will load the u.item.txt file so that we can extract records from the basic.csv dataset for the movies in the u.item.txt file.

We will use the DataFrame isin() function to find all the movies from the imdb movies file. Once we have identified all the movies, we will loop through each movie and use its ID (tconst) to find its rating and total vote count. We will then calculate the weighted average rating for each movie in the same way we did in Q5.

Implementation of the current process is shared below,

```python
1  import pandas as pd
2  from pandas import DataFrame
3  import numpy as np
4
5  load_movie_data = pd.read_csv('basic.tsv', delimiter="\t")
6  load_rating_data = pd.read_csv('rating.tsv', delimiter="\t")
7
8  col = ['movie_id', 'movie_title', 'release_date', 'video_release_date',
         'IMDb_URL', 'unknown', 'Action', 'Adventure', 'Animation', '
      Childrens',
9          'Comedy', 'Crime', 'Documentary', 'Drama', 'Fantasy', 'Film_Noir
      ', 'Horror', 'Musical', 'Mystery', 'Romance', 'Sci_Fi', 'Thriller',
      'War', 'Western']
10 item_dataset = pd.read_csv('u.item.txt', delimiter="|", header=None,
      names=col)
11
12 list_movies = []
13
14 list_titles = item_dataset.iloc[:, 1]
15
16 filtered_movies = load_movie_data[load_movie_data['primaryTitle'].isin(
      list_titles)]
17
18 print(filtered_movies)
19
20 list_rating = []
21
22 for movie in filtered_movies.iterrows():
23     filtered_rating = load_rating_data[load_rating_data['tconst'] ==
```

```
        movie[1]['tconst']]
24      if len(filtered_rating.index) > 0:
25          avg_rating = filtered_rating.iloc[0,1]
26          num_Votes = filtered_rating.iloc[0,2]
27          rating_item = {
28              'tconst': movie[1]['tconst'],
29              'movie_title': movie[1]['primaryTitle'],
30              'avg_rating': avg_rating,
31              'num_votes': num_Votes,
32              'weighted_rating':  avg_rating * num_Votes
33          }
34          list_rating.append(rating_item)
35
36 final_rating_dataset = DataFrame(list_rating)
37 final_rating_dataset = final_rating_dataset.sort_values('
       weighted_rating', ascending=False)
38 print(final_rating_dataset.iloc[0:10,:])
39 print(final_rating_dataset.iloc[-10:])
```

**Listing 6:** Ranking the movies

After calculating the weighted average rating of each movie, we sorted the dataframe based on this column and then used the iloc function to get the top and bottom 10 movies. Finally, we saved the complete ranked result in a separate file. The output of the top and bottom 10 movies is shown below.

**Table 15:** Top 10 ranked movies

| Movies | Avg. Rating | Votes | Weighted Average |
|--------|-------------|-------|------------------|
| Batman (1989) | 7.6 | 64 | 486.4 |
| Cinderella (1950) | 7.6 | 63 | 478.8 |
| Beauty and the Beast (1991) | 7.2 | 57 | 410.4 |
| Scream 2 (1997) | 8.4 | 21 | 176.4 |
| Wes Craven's New Nightmare (1994) | 8.8 | 20 | 176.0 |
| Scream (1996) | 8.5 | 19 | 161.5 |
| Alien (1979) | 7.0 | 22 | 154.0 |
| Night of the Living Dead (1968) | 8.4 | 18 | 151.2 |
| Halloween: The Curse of Michael Myers (1995) | 8.3 | 18 | 149.4 |
| Jurassic Park (1993) | 8.6 | 17 | 146.2 |

**Table 16:** Least 10 ranked movies

| Movies | Avg. Rating | Count | Weighted Average |
|---|---|---|---|
| Akira (1988) | 7.2 | 11 | 79.2 |
| Super Mario Bros. (1993) | 7.4 | 9 | 66.6 |
| Spawn (1997) | 7.1 | 9 | 63.9 |
| Deep Rising (1998) | 8.8 | 6 | 52.8 |
| Blade Runner (1982) | 7.4 | 7 | 51.8 |
| Top Gun (1986) | 7.3 | 7 | 51.1 |
| Star Trek VI: The Undiscovered Country (1991) | 7.7 | 6 | 46.2 |
| Rumble in the Bronx (1995) | 8.6 | 5 | 43.0 |
| Return of the Jedi (1983) | 8.6 | 5 | 43.0 |
| Bride of Frankenstein (1935) | 4.2 | 5 | 21.0 |

# References

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

- DataFrame operations, `https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.groupby.html`

- Tweepy, `https://docs.tweepy.org/en/latest/`

- Numpy, `https://numpy.org/doc/stable/reference/generated/numpy.log2.html`