

ex45

August 12, 2022

```
[25]: # Solution Ex. 45
```

```
[3]: #inputPathWatched = "/data/students/bigdata-01QYD/ex_data/Ex45/data/  
    ↪watchedmovies.txt"  
#inputPathPreferences = "/data/students/bigdata-01QYD/ex_data/Ex45/data/  
    ↪preferences.txt"  
#inputPathMovies = "/data/students/bigdata-01QYD/ex_data/Ex45/data/movies.txt"  
#outputPath = "res_out_Ex45/"  
#threshold = 0.5  
  
inputPathWatched = "data/Ex45/data/watchedmovies.txt"  
inputPathPreferences = "data/Ex45/data/preferences.txt"  
inputPathMovies = "data/Ex45/data/movies.txt"  
outputPath = "res_out_Ex45/"  
threshold = 0.5
```

```
[4]: # Read the content of the watched movies file  
watchedRDD = sc.textFile(inputPathWatched)
```

```
[5]: # Select only userid and movieid  
# Define an RDD of pairs with movieid as key and userid as value  
movieUserPairRDD = watchedRDD.map(lambda line: (line.split(",")[1], line.  
    ↪split(",")[0]))
```

```
[6]: # Read the content of the movies file  
moviesRDD = sc.textFile(inputPathMovies)
```

```
[7]: # Select only movieid and genre  
# Define an RDD of pairs with movieid as key and genre as value  
movieGenrePairRDD = moviesRDD.map(lambda line: (line.split(",")[0], line.  
    ↪split(",")[2]))
```

```
[8]: # Join watched movies with movies  
joinWatchedGenreRDD = movieUserPairRDD.join(movieGenrePairRDD)
```

```
[9]: # Select only userid (as key) and genre (as value)
```

```
usersWatchedGenresRDD = joinWatchedGenreRDD.map(lambda pair: (pair[1][0],  
↪pair[1][1]))
```

```
[10]: # Read the content of preferences.txt  
preferencesRDD = sc.textFile(inputPathPreferences)
```

```
[11]: # Define an RDD of pairs with userid as key and genre as value  
userLikedGenresRDD = preferencesRDD.map(lambda line: (line.split(",")[0], line.  
↪split(",")[1]))
```

```
[12]: # Cogroup the lists of watched and liked genres for each user  
# There is one pair for each userid  
# the value contains the list of genres (with repetitions) of the  
# watched movies and the list of liked genres  
userWatchedLikedGenres = usersWatchedGenresRDD.cogroup(userLikedGenresRDD)
```

```
[13]: # This function is used in the next transformation to select users with a  
↪misleading profile  
def misleadingProfileFunc(userWatchedLikedGenresLists):  
    # Store in a local list the "small" set of liked genres  
    # associated with the current user  
    likedGenres = list(userWatchedLikedGenresLists[1][1])  
  
    # Iterate over the watched movies (the genres of the watched movies) and  
↪count  
    # - The number of watched movies for this user  
    # - How many of watched movies are associated with a not liked genre  
    numWatchedMovies = 0  
    notLiked = 0  
  
    for watchedGenre in userWatchedLikedGenresLists[1][0]:  
        numWatchedMovies = numWatchedMovies+1  
        if watchedGenre not in likedGenres:  
            notLiked = notLiked+1  
  
    # Check if the number of watched movies associated with a non-liked genre  
    # is greater than threshold%  
    if float(notLiked) > threshold * float(numWatchedMovies):  
        return True  
    else:  
        return False
```

```
[14]: # Filter the users with a misleading profile  
misleadingUsersListsRDD = userWatchedLikedGenres.filter(misleadingProfileFunc)
```

```
[15]: # This function is used in the next transformation to select the pairs
      ↪(userid,misleading genre)
def misleadingGenresFunc(userWatchedLikedGenresLists):
    # Store in a local list the "small" set of liked genres
    # associated with the current user

    userId = userWatchedLikedGenresLists[0]
    likedGenres = list(userWatchedLikedGenresLists[1][1])

    # In this solution I suppose that the number of distinct genres for each
    ↪user
    # is small and can be stored in a local variable.
    # The local variable is a dictionary that stores for each non-liked genre
    # also its number of occurrences in the list of watched movies of the
    ↪current user
    numGenres = {}

    # Iterate over the watched movies (the genres of the watched movies).
    # Select the watched genres that are not in the liked genres and
    # count their number of occurrences. Store them in the numGenres dictionary
    for watchedGenre in userWatchedLikedGenresLists[1][0]:
        # Check if the genre is not in the liked ones
        if watchedGenre not in likedGenres:
            # Update the number of times this genre appears
            # in the list of movies watched by the current user
            if watchedGenre in numGenres:
                numGenres[watchedGenre] = numGenres[watchedGenre] + 1
            else:
                numGenres[watchedGenre] = 1

    # Select the genres, which are not in the liked ones,
    # which occur at least 5 times
    selectedGenres = []

    for genre, occurrences in numGenres.items():
        if occurrences>=5:
            selectedGenres.append( (userId, genre) )

    return selectedGenres
```

```
[16]: # Select the pairs (userid,misleading genre)
misleadingUserGenrePairRDD = misleadingUsersListsRDD.
      ↪flatMap(misleadingGenresFunc)
```

```
[18]: #misleadingUserGenrePairRDD.collect()
```

```
[ ]: misleadingUserGenrePairRDD.saveAsTextFile(outputPath)
```