

Data Science [ITE4005]

Programming Assignment #4

: Predict the ratings in test data using training data

2015004120

Software Engineering

Park Seonha

0. Environment

- OS : Windows
- Language : Python3 (version : 3.4.4)
- Executable file : recommender.py
(../Programming_Assignment_4/project_recommender/recommender.py)

1. Summary of Algorithm

This programming assignment's goal is "Predict the ratings of movies in test data by using the given training data containing movie ratings of users" with any algorithm which can use to predict such as content-based algorithm and collaborative filtering algorithms. I used user based nearest neighbor collaborative filtering algorithm to predict the ratings.

Before predict the rating, I calculated all necessary data like similarity and neighbor set because of its hard computation amount.

* precomputed data

- mean of all each users
- list of neighbor of each user
- similarity between all users.

To calculate similarity I used below formula.

$$sim(x, y) = \frac{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_s)(r_{ys} - \bar{r}_y)}{\sqrt{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_s)^2} \sqrt{\sum_{s \in S_{xy}} (r_{ys} - \bar{r}_y)^2}}$$

After calculate all necessary data, calculate the predict value with user_id and item_id. find the neighbors who have the rating data of certain item_id, and get the similarity between the user_id and each neighbor. To predict the user a's rating of item p, I used this formula.

$$pred(a, p) = \bar{r}_a + \frac{\sum_{b \in N} sim(a, b) * (r_{b,p} - \bar{r}_b)}{\sum_{b \in N} sim(a, b)}$$

And check the boundary of rating to fit in result value.

2. Description of codes

```
77 def main():
78     # get command line argument
79     train_file = sys.argv[1]
80     test_file = sys.argv[2]
```

get training file name and test file name from command line argument.

```
82 with open(train_file) as f:
83     train_data = f.readlines()
84     train_data = [d.strip() for d in train_data]
85
86     # make record_list
87     record_list = {}
88     user_list = []
89     item_list = []
90     user_list = set(user_list)
91     item_list = set(item_list)
92     for line in train_data:
93         each_line = line.split("\t")
94         temp_key = (int(each_line[0]), int(each_line[1]))
95         temp_value = int(each_line[2])
96         record_list[temp_key] = temp_value
97         user_list.add(int(each_line[0]))
98         item_list.add(int(each_line[1]))
99
```

parse training file into record_list(dictionary whose key is tuple of user and item and whose value is the rating), user_list (set of all user_ids exist in training data), item_list(set of all item_id exist in training data)

```
100 # make user_item
101 user_item = {}
102 user_mean = {}
103 for user in user_list:
104     mean = 0.0
105     for item in item_list:
106         if (user,item) in record_list:
107             if user in user_item.keys():
108                 user_item[user].append(item)
109                 mean += record_list[(user,item)]
110             else:
111                 user_item[user] = []
112                 user_item[user].append(item)
113     mean = float(mean) / float(len(user_item[user]))
114     user_mean[user] = mean
```

pre-compute all the mean of each users and make the dictionary whose key is user_id and value is the list of its item at the same time.

```
118 # make get_intersect and get common average
119 nume = 0
120 deno = 0
121 get_intersect = {}
122 for k1, v1 in user_item.items():
123     for k2, v2 in user_item.items():
124         if k1 != k2:
125             intersected = set(user_item[k1]).intersection((user_item[k2]))
126             get_intersect[(k1,k2)] = intersected
127             nume += len(intersected)
128             deno += 1
129
130 average_common = int(float(nume)/float(deno))
131
132 # make sim dict
133 sim_dict = {}
134 for k,v in get_intersect.items():
135     if len(v) > average_common:
136         sim_dict[k] = get_sim(k[0],k[1],v,record_list)
```

make dictionary whose key is the tuple of two user and value is the set which is

intersection of two user's item. And pre-compute the average of it at same time.

pre-compute all similarities in get_intersect dictionary and store them in sim_dict with same key.

```

138     # make neighbor dict
139     neighbors = {}
140     for user1 in user_list:
141         for user2 in user_list:
142             if user1 != user2 :
143                 if len(get_intersect[(user1,user2)]) > average_common:
144                     if user1 in neighbors.keys():
145                         neighbors[user1].append(user2)
146                     else:
147                         neighbors[user1] = []
148                         neighbors[user1].append(user2)
149             if user1 not in neighbors.keys():
150                 neighbors[user1] = []

```

precompute all user's neighbor(the intersection of both user's item is larger than average common items) and store it into dictionary whose key is user_id and value is the list of neighbor users.

```

8  def get_sim(user1, user2, intersected, record_list):
9      mean_user1 = 0.0
10     for item in intersected:
11         mean_user1 += record_list[(user1,item)]
12     mean_user1 = float(mean_user1)/float(len(intersected))
13
14     mean_user2 = 0.0
15     for item in intersected:
16         mean_user2 += record_list[(user2,item)]
17     mean_user2 = float(mean_user2)/float(len(intersected))
18
19     dist_user1 = 0.0
20     dist_list1 = []
21     for item in intersected:
22         temp = float(record_list[(user1,item)])-mean_user1
23         dist_list1.append(temp)
24         dist_user1 += math.pow(temp,2.0)
25     dist_user1 = math.sqrt(dist_user1)
26
27     dist_user2 = 0.0
28     dist_list2 = []
29     for item in intersected:
30         temp = float(record_list[(user2,item)])-mean_user2
31         dist_list2.append(temp)
32         dist_user2 += math.pow(temp,2.0)
33     dist_user2 = math.sqrt(dist_user2)
34
35     sim_nume = 0.0
36     for i,j in zip(dist_list1,dist_list2):
37         sim_nume += i*j
38
39     sim_deno = dist_user1 * dist_user2
40     if abs(sim_deno) < 0.00001:
41         return 1.0
42
43
44     return sim_nume/sim_deno

```

compute the formula
$$sim(x,y) = \frac{\sum_{seSxy} (r_{xs} - \bar{r}_s)(r_{ys} - \bar{r}_s)}{\sqrt{\sum_{seSxy} (r_{xs} - \bar{r}_s)^2} \sqrt{\sum_{seSxy} (r_{ys} - \bar{r}_s)^2}}$$
 between two users.

```

152     #start test
153     with open(test_file) as f:
154         test_data = f.readlines()
155         test_data = [d.strip() for d in test_data]
156
157     test_list = []
158     for line in test_data:
159         each_line = line.split("\t")
160         temp_tuple = (int(each_line[0]), int(each_line[1]))
161         test_list.append(temp_tuple)
162

```

open test data file and make the list of tuple (test_user_id, test_item) which to predict.

```

164     output_file = train_file+"_prediction.txt"
165     out_f = open(output_file,"w")
166     for test in test_list:
167         out_f.write(str(test[0])+"\t"+str(test[1])+"\t")
168         out_f.write(str(predict(test[0],test[1],record_list,user_item,get_intersect,average_common,user_mean,sim_dict,neighbors
169         out_f.write("\n")
170     out_f.close()
171

```

predict the rating of each test and write in ouput file.

```

47 def predict(user, item, record_list, user_item, get_intersect,average_common, user_mean,sim_dict,neighbor_dict):
48     mean_user = user_mean[user]
49
50     neighbor = []
51     for i in neighbor_dict[user]:
52         if item in user_item[i]:
53             neighbor.append(i)
54
55     sim_list = []
56     sim_l_num = []
57
58     for neigh in neighbor:
59         temp = sim_dict[(user,neigh)]
60         sim_list.append(temp)
61         sim_l_num.append(temp * (record_list[(neigh,item)] - user_mean[neigh]))
62
63     if len(sim_list) > 0 and sum(sim_list) > 0 :
64         result = mean_user + (sum(sim_l_num)/sum(sim_list))
65     else:
66         result = mean_user
67
68     if result < 1:
69         result = 1
70     elif result > 5:
71         result = 5
72     else:
73         result = round(result)
74
75     return result

```

predict the rating of *item* which *user* rated with using pre computed datas and the

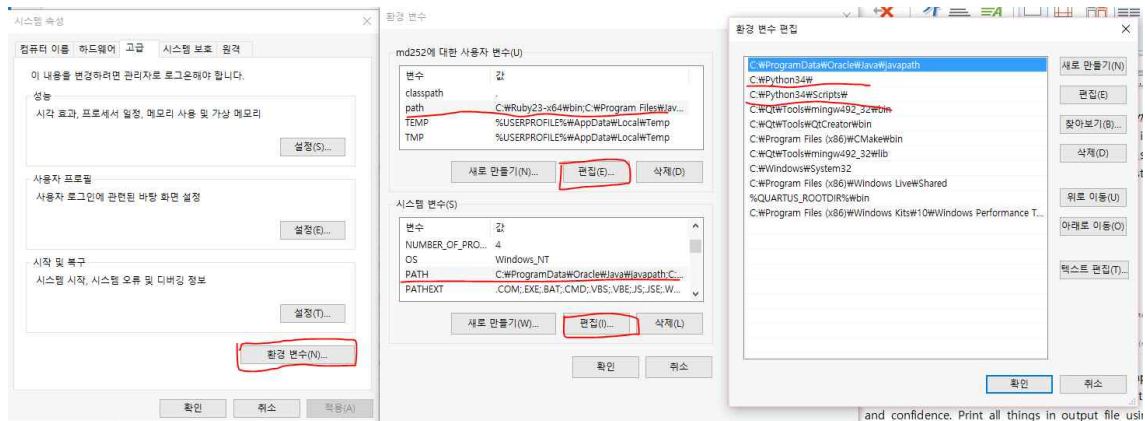
formula
$$pred(a,p) = \bar{r}_a + \frac{\sum_{b \in N} sim(a,b) * (r_{b,p} - \bar{r}_b)}{\sum_{b \in N} sim(a,b)}$$
 . If predict value is larger than 5, return 5. And if predict value is smaller than 1, return 1. Otherwise, return the round value of predicted result.

3. Instructions for compiling this code

This code is written in Python3 and tested in Python 3.4.4. Since this code has **sqrt** and **pow**(to calculate similarity between two user), and it exists in Python3 math module, to run this code, we need python3 interpreter.

We can get python 3.4.4 in (<https://www.python.org/downloads/release/python-344/>)

If python2 is already installed in PC, change path Python2x to Python34 in advanced system settings – environment variables.

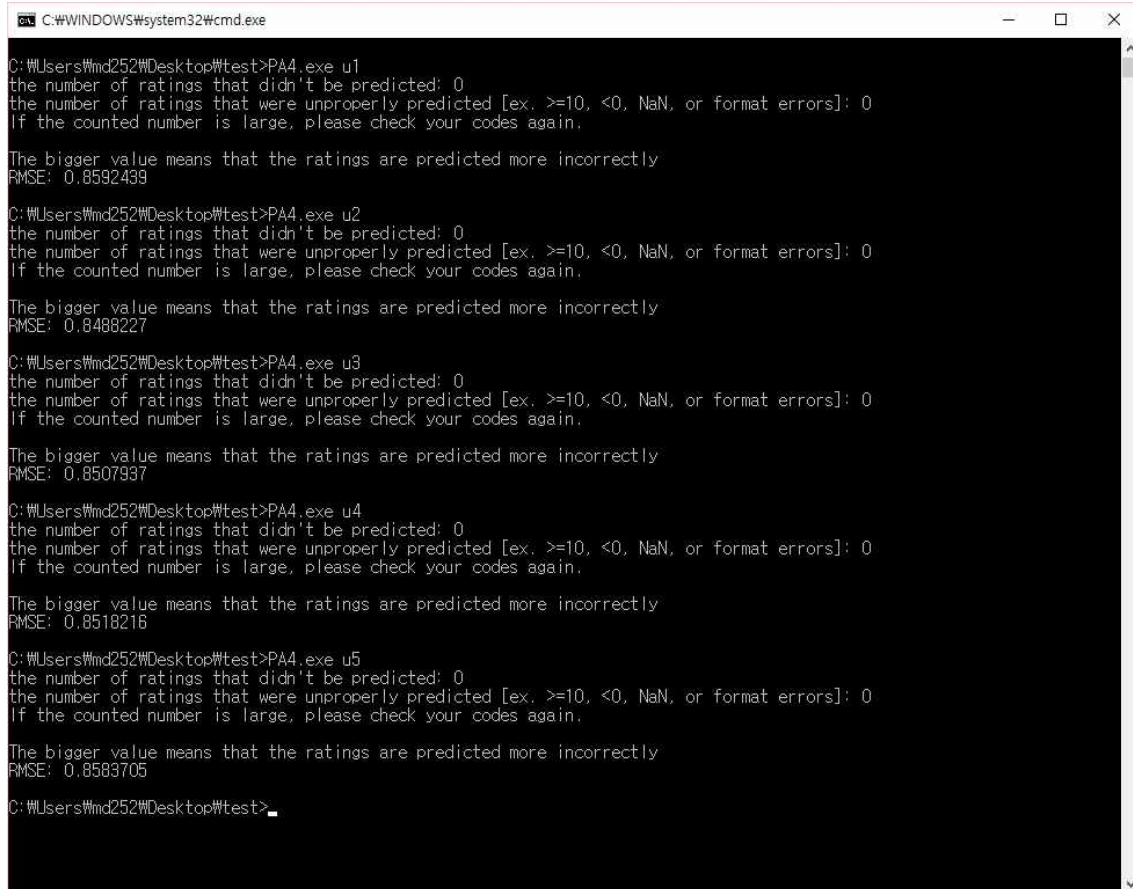


After environment variables setting, we can get python3 default in cmd.

In code directory(python file is in /project_recommender), we can run this program on cmd "python recommender.py [training file name] [test file name]".

The output file will be located in same directory with training file.

4. Other Specifications



```
C:\WINDOWS\system32\cmd.exe
C:\Users\Wmd252\Desktop\#test>PA4.exe u1
the number of ratings that didn't be predicted: 0
the number of ratings that were improperly predicted [ex. >=10, <0, NaN, or format errors]: 0
If the counted number is large, please check your codes again.

The bigger value means that the ratings are predicted more incorrectly
RMSE: 0.8592439

C:\Users\Wmd252\Desktop\#test>PA4.exe u2
the number of ratings that didn't be predicted: 0
the number of ratings that were improperly predicted [ex. >=10, <0, NaN, or format errors]: 0
If the counted number is large, please check your codes again.

The bigger value means that the ratings are predicted more incorrectly
RMSE: 0.8498227

C:\Users\Wmd252\Desktop\#test>PA4.exe u3
the number of ratings that didn't be predicted: 0
the number of ratings that were improperly predicted [ex. >=10, <0, NaN, or format errors]: 0
If the counted number is large, please check your codes again.

The bigger value means that the ratings are predicted more incorrectly
RMSE: 0.8507937

C:\Users\Wmd252\Desktop\#test>PA4.exe u4
the number of ratings that didn't be predicted: 0
the number of ratings that were improperly predicted [ex. >=10, <0, NaN, or format errors]: 0
If the counted number is large, please check your codes again.

The bigger value means that the ratings are predicted more incorrectly
RMSE: 0.8518216

C:\Users\Wmd252\Desktop\#test>PA4.exe u5
the number of ratings that didn't be predicted: 0
the number of ratings that were improperly predicted [ex. >=10, <0, NaN, or format errors]: 0
If the counted number is large, please check your codes again.

The bigger value means that the ratings are predicted more incorrectly
RMSE: 0.8583705

C:\Users\Wmd252\Desktop\#test>_
```

In testing, some computer can take quite long time to pre-compute similarity before predict. I only calculate the similarity between each neighbor just one time, it is necessary step and it can't make infinite loop, wait until it will finish.