

EECS E6893 Big Data Analytic HW2

Chong Hu ch3467

October 18, 2019

Problem 1. Friends Recommendation

Result:

19/10/11 01:12:02 WARN Nati Using Spark's default log4j Setting default log level to To adjust logging level use For user 13420: (4736, 3) (7651, 3) (10469, 3) (14264, 3) (351, 2) (2101, 2) (2554, 2) (7608, 2) (8508, 2) (8711, 2) For user 8942: (8939, 3) (8940, 1) (8943, 1) (8944, 1) For user 5850: (5819, 3) (5805, 2) (5811, 2) (5815, 2) (5828, 2) (5831, 2) (5836, 2) (219, 1) (576, 1) (639, 1)	For user 49824: (49846, 3) (41581, 2) (43382, 2) (49786, 2) (49788, 2) (49789, 2) (49814, 2) (49819, 2) (49834, 2) (16, 1) For user 924: (439, 1) (2409, 1) (6995, 1) (11860, 1) (15416, 1) (43748, 1) (45881, 1) For user 44410: (4231, 3) (44462, 3) (351, 2) (4302, 2) (6318, 2) (8221, 2) (9095, 2) (10328, 2) (10370, 2) (10462, 2)	For user 8974: (8960, 12) (12241, 12) (8774, 10) (6973, 7) (8969, 6) (8980, 4) (8982, 4) (8984, 4) (8978, 3) (8979, 3) For user 8941: (8943, 2) (8944, 2) (8940, 1) For user 9019: (9022, 2) (317, 1) (9023, 1) For user 9993: (9991, 5) (13134, 1) (13478, 1) (13877, 1) (34299, 1) (34485, 1) (34642, 1) (37941, 1)
--	---	--

Codes:

```
1 from pyspark import SparkConf, SparkContext
2 import pyspark
3 import sys
4
5
6 def getData(sc, filename):
7     """
8     Load data from raw text file into RDD and transform.
```

```

9      Hint: transformation you will use: map(<lambda function>).
10     Args:
11         sc (SparkContext): spark context.
12         filename (string): hw2.txt cloud storage URI.
13     Returns:
14         RDD: RDD list of tuple of (<User>, [friend1, friend2, ... ]),
15         each user and a list of user's friends
16     """
17     # read text file into RDD
18     data = sc.textFile(filename)
19     data = data.map(lambda line: line.split("\t")).map(
20         lambda line: (int(line[0]), [int(x) for x in line[1].split(",")] if len(
21             line[1]) else []))
22     return data
23
24
25 def map_friends(line):
26     """
27     map function to construct edge between friends
28     construct a pair of ((friend1, friend2) -> common friends list)
29     if two friends are already direct friends, then common friends list
30     is empty.
31     :param line: tuple of (<User>, [friend1, friend2, ... ]),
32     each user and a list of user's friends
33     :return:
34     """
35     user = line[0]
36     friends = line[1]
37     yield ((user, user), [])
38     for i in range(len(friends)):
39         yield ((user, friends[i]), [])
40         for j in range(len(friends)):
41             yield ((friends[i], friends[j]), [user])
42             yield ((friends[j], friends[i]), [user])
43
44
45 def reduce_friends_pair(pair1, pair2):
46     """
47     reduce function to reduce same friend-friend pairs.
48     If the common friend list is None, which means they are direct friends,
49     still return empty list
50     :param pair1: reduceByKey first pair
51     :param pair2: reduceByKey second pair
52     :return: a pair of ((friend1, friend2) -> common friends list)
53     if two friends are already direct friends,

```

```

54         then common friends list is empty.
55     """
56     if len(pair1) == 0 or len(pair2) == 0:
57         return []
58     common_friends = set(pair1).union(set(pair2))
59     return list(common_friends)
60
61
62 def find_mutual(line):
63     """
64     map mutual edges into (user, (friend, common_friend_num))
65     :param line: a pair of ((friend1, friend2) -> common friends list)
66                 if two friends are already direct friends,
67                 then common friends list is empty.
68     :return: (user, (friend, common_friend_num))
69     """
70     return line[0][0], (line[0][1], len(line[1]))
71
72
73 def sort_top_friends(line):
74     """
75     sort friend has most common friends
76     :param line: (user, [(friend, common_friend_num), ...])
77     :return: recommendations result for the user
78     """
79     user = line[0]
80     friends = sorted(line[1], key=lambda x: (x[1], -x[0]), reverse=True)
81     while len(friends) > 0:
82         if friends[-1][1] == 0:
83             friends.pop()
84         else:
85             break
86     if len(friends) > 10:
87         friends = friends[0:10]
88     return user, friends
89
90
91 def main():
92     # Configure Spark
93     sc = pyspark.SparkContext.getOrCreate()
94     # The directory for the file
95     filename = "q1.txt"
96
97     # Get data in proper format
98     data = getData(sc, filename)

```

```

99
100     # Get set of all mutual friends
101     mapData = data.flatMap(map_friends).reduceByKey(reduce_friends_pair)
102     # print(mapData.take(10))
103     # For each person, get top 10 mutual friends
104     getFriends = mapData.map(find_mutual).groupByKey().map(sort_top_friends)
105     # print(getFriends.take(5))
106     # Only save the ones we want
107     wanted = [924, 8941, 8942, 9019, 49824, 13420, 44410, 8974, 5850, 9993]
108     result = getFriends.filter(lambda x: x[0] in wanted).collect()
109     for res in result:
110         print("For user %d:" % res[0])
111         for recommendation in res[1]:
112             print(recommendation)
113         print()
114     sc.stop()
115
116
117 if __name__ == "__main__":
118     main()

```

Problem 2. Graph Analysis

Connected Components:

```

clusters amount: 917

number of users in top 10 cluster
cluster id: 0      number of users: 48860
cluster id: 38403  number of users: 66
cluster id: 18466  number of users: 31
cluster id: 18233  number of users: 25
cluster id: 18891  number of users: 19
cluster id: 864    number of users: 16
cluster id: 49297  number of users: 13
cluster id: 19199  number of users: 6
cluster id: 7658   number of users: 5
cluster id: 22897  number of users: 4
Total number of users in top 10 cluster: 49045

user ids for the cluster which has 25 users
[18233, 18234, 18235, 18236, 18237, 18238, 18239, 18240, 18241, 18242, 18243, 18244, 18245, 18246, 18247, 18248, 18249, 18250, 18251, 18252, 18253, 18254, 18255, 18256, 18257]

```

(1) How many clusters / connected components in total for this dataset?

```
1 clusters amount: 917
```

(2) How many users in the top 10 clusters? There are different number of users in each clusters, so rank them and give the top 10 clusters with the largest amount of users.

```
1 number of users in top 10 cluster
2 cluster id:      0      number of users: 48860
3 cluster id:      38403   number of users: 66
```

```

4 cluster id:      18466      number of users:      31
5 cluster id:      18233      number of users:      25
6 cluster id:      18891      number of users:      19
7 cluster id:      864        number of users:      16
8 cluster id:      49297      number of users:      13
9 cluster id:      19199      number of users:      6
10 cluster id:      7658       number of users:      5
11 cluster id:      22897      number of users:      4
12 Total number of users in top 10 cluster:      49045

```

- (3) What are the user ids for the cluster which has 25 users? Basically, list out all the 25 user IDs in that cluster.

```

1 user ids for the cluster which has 25 users
2 [18233, 18234, 18235, 18236, 18237, 18238, 18239, 18240, 18241, 18242, 18243, 18244,
3 18245, 18246, 18247, 18248, 18249, 18250, 18251, 18252, 18253, 18254, 18255, 18256,
4 18257]

```

Page rank

```

a list of 10 important users (User ID) in this network:
[10164, 15496, 14689, 24966, 7884, 934, 45870, 5148, 20283, 46039]
The most important one is 10164

```

```

Using different parameter:
a list of 10 important users (User ID) in this network:
[10164, 15496, 14689, 24966, 7884, 934, 45870, 20283, 46039, 14996]
The most important one is 10164

```

- (4) Provide a list of 10 important users (User ID) in this network. Who is the most important one? Order by the “PageRank” value.

```

1 a list of 10 important users (User ID) in this network:
2 [10164, 15496, 14689, 24966, 7884, 934, 45870, 5148, 20283, 46039]
3 The most important one is 10164

```

- (5) By using different parameter settings for PageRank, is there any difference? This is an open question, you can try as many as you want. Provide the screenshots of your tests.

```

1 Using different parameter:
2 a list of 10 important users (User ID) in this network:
3 [10164, 15496, 14689, 24966, 7884, 934, 45870, 20283, 46039, 14996]
4 The most important one is 10164

```

Although we used different parameter settings for PageRank, we can still get similar result. We can see many common users in both result and the most important one is the same.

- 6 Why this user become the most important one? What are the possible reasons? This is an open question, basically, understand how PageRank works. You can also use the result from the connected component to explain it.

From the connected component analysis, I can find this user is in the largest cluster. This makes this user relatively important. Beyond that, most of this user's friend might have high PageRank value. This user's friends who have high PageRank value contributes to this user.

7 PageRank Calculation

	ID1	ID2	ID3	ID4	ID5
Iteration 0	0.2	0.2	0.2	0.2	0.2
Iteration 1	0.07	0.29	0.41	0.07	0.16
Iteration 2	0.09	0.45	0.25	0.09	0.12
Iteration 3	0.13	0.29	0.29	0.13	0.16
Iteration 4	0.09	0.34	0.33	0.09	0.15

Table 2: PageRank Calculation Details

Here is whole code for part II.

```
1 import pyspark
2 from pyspark import SparkConf, SparkContext
3 from pyspark import SQLContext
4 import os
5 from graphframes import *
6
7
8 def getData(sc, filename):
9     """
10     Load data from raw text file into RDD and transform.
11     Hint: transformation you will use: map(<lambda function>).
12     Args:
13         sc (SparkContext): spark context.
14         filename (string): hw2.txt cloud storage URI.
15     Returns:
16         RDD: RDD list of tuple of (<User>, [friend1, friend2, ... ]),
17         each user and a list of user's friends
18     """
```

```

19     # read text file into RDD
20     data = sc.textFile(filename)
21     data = data.map(lambda line: line.split("\t")).map(
22         lambda line: (int(line[0]), [int(x) for x in line[1].split(",")] if len(
23             line[1]) else []))
24     return data
25
26
27 def get_vertices(data, sqlcontext):
28     """
29     get vertices
30     :param data: RDD list of tuple of (<User>, [friend1, friend2, ... ]),
31     each user and a list of user's friends
32     :param sqlcontext: SQLContext
33     :return: dataframe
34     """
35     vertices = data.map(lambda line: (line[0],))
36
37     return sqlcontext.createDataFrame(vertices, schema=["id"])
38
39
40 def get_edges(data, sqlcontext):
41     """
42     get edges
43     :param data: RDD list of tuple of (<User>, [friend1, friend2, ... ]),
44     each user and a list of user's friends
45     :param sqlcontext: SQLContext
46     :return:
47     """
48
49 def map_friends(line):
50     """
51     map function to construct edge between friends
52     construct a pair of ((friend1, friend2) -> common friends list)
53     if two friends are already direct friends, then common friends list
54     is empty.
55     :param line: tuple of (<User>, [friend1, friend2, ... ]),
56     each user and a list of user's friends
57     :return: friend pair
58     """
59     user = line[0]
60     friends = line[1]
61     for i in range(len(friends)):
62         yield (user, friends[i])
63

```

```

64     edges = data.flatMap(map_friends)
65     return sqlcontext.createDataFrame(edges, schema=["src", "dst"])
66
67
68 def connected_components(graph):
69     """
70     run connected components on graph
71     :param graph: Graph contains vertices and edges
72     :return:
73     """
74     print("connected components")
75     result = graph.connectedComponents()
76
77     # How many clusters / connected components in total for this dataset
78     cluster_num = result.select("component").distinct().count()
79     print("clusters amount: ", cluster_num)
80     print()
81
82     # How many users in the top 10 clusters?
83     print("number of users in top 10 cluster")
84     res1 = result.groupBy("component").count().orderBy('count',
85                                                         ascending=False)
86
87     res2 = res1.head(10)
88     total = 0
89     for row in res2:
90         total += row["count"]
91         print("cluster id:\t%d\tnumber of users:\t%d" % (
92             row["component"], row["count"]))
93     print("Total number of users in top 10 cluster:\t", total)
94     print()
95
96     # What are the user ids for the cluster which has 25 users?
97     print("user ids for the cluster which has 25 users")
98     cluster_id = res1.where(res1["count"] == 25).select("component").collect()
99     cluster_id = [row["component"] for row in cluster_id]
100     user_list = result.where(result["component"].isin(cluster_id)).select(
101         "id").collect()
102     user_ls = [row["id"] for row in user_list]
103     print(user_ls)
104     print()
105     return
106
107 def page_rank(graph):
108     """

```



```

109     run PageRank on graph
110     :param graph: Graph contains vertices and edges
111     :return:
112     """
113     print("PageRank:")
114     result = graph.pageRank(resetProbability=0.15, tol=0.01)
115
116     # Provide a list of 10 important users (User ID) in this network.
117     print("a list of 10 important users (User ID) in this network:")
118     user_list = result.vertices.select("id", "pagerank") \
119         .orderBy('pagerank',
120             ascending=False).head(10)
121     user_ls = [row["id"] for row in user_list]
122     print(user_ls)
123     print("The most important one is %d" % user_ls[0])
124     print()
125
126     # using different parameter settings for PageRank
127     print("Using different parameter:")
128     result = graph.pageRank(resetProbability=0.1, maxIter=20)
129     print("a list of 10 important users (User ID) in this network:")
130     user_list = result.vertices.select("id", "pagerank") \
131         .orderBy('pagerank',
132             ascending=False).head(10)
133     user_ls = [row["id"] for row in user_list]
134     print(user_ls)
135     print("The most important one is %d" % user_ls[0])
136
137
138 def main():
139     # Configure Spark
140     if not os.path.isdir("checkpoints"):
141         os.mkdir("checkpoints")
142     conf = SparkConf().setMaster('local').setAppName('connected components')
143     sc = SparkContext(conf=conf)
144     sqlcontext = SQLContext(sc)
145     SparkContext.setCheckpointDir(sc, "checkpoints")
146
147     # The directory for the file
148     filename = "../q1/q1.txt"
149
150     # Get data in proper format
151     data = getData(sc, filename)
152     edges = get_edges(data, sqlcontext)
153     vertices = get_vertices(data, sqlcontext)

```

```
154     graph = GraphFrame(vertices, edges)
155     connected_components(graph=graph)
156     page_rank(graph=graph)
157
158
159 if __name__ == '__main__':
160     main()
```
