

CS838 PROJECT REPORT - STAGE 4

COMBINING TWO TABLES

(04/ 16/ 2017)

Gaurav Mishra
gmishra2@wisc.edu

Om Jadhav
ojadhav@wisc.edu

Pallavi M Kakunje
kakunje@wisc.edu

1. Overview:

The goal of this project stage is to combine the matching tuples i.e., tuples that represent same entity, from two tables into a single table. In this stage, we used the same data sources as in the previous stages.

Entity type: Mobile game application

Data Sources: Google Play Store (google table) and Apple Store (apple table).

We use the matches obtained from Stage 3 of the project and combine the corresponding matched tuples in this stage.

2. Implementation:

2.1. Combining Process:

Our best matcher as identified from the previous stage (Stage 3) is Random Forest. This matcher is applied to all the candidate matches that survived the blocking stage. This provided all the matches. Resulting matches are stored in all_predictions.csv.

Our python script goes through all the predictions to find whether the tuples are predicted as match or not. If there is a match, then we copy corresponding tuples from both the tables into matched_tuples.csv. Further we merge both tuples according to merger rules and copy the merged tuple into merged_tuples.csv (table E).

This stage involved combining just the two above mentioned tables and no additional table. The process of combining was smooth, without facing any unexpected issues.

2.2 Merging rules and Schema

Both google table and apple table have the same schema of the form:

(*"ID"*, *"name"*, *"category"*, *"developer"*, *"rating"*)

The final merged table also has the same schema. The attribute values are merged using following rules: ID takes new values in the merged table E. "name" and "developer" are same in both table if there is a match. So, we copy it from one of the table (google table is the chosen one). After examining the tables, we observed that "category" in google table makes more sense as compared to "category" in apple table. So we pick category from google table. "Rating" is computed as the average value of "rating" values from both the tables. There are few tuples in which rating is not present in apple table, in this case we pick rating from google table alone.

2.3. Merged Table (Table E):

Both google table and apple table have same schema. So, merged_tuples table (table E) will also have same schema as mentioned above. Each tuple of E contains attributes: "ID", "name", "category", "developer" and "rating".

There are total 490 tuples in table E.

Few tuples of table E are as follows:

ID	name	category	developer	rating
1	Candy Crush Saga	Casual	King	4.554035
2	Super Mario Run	Action	Nintendo Co., Ltd.	3.33734
3	Jetpack Joyride	Arcade	Halfbrick Studios	4.421695
4	Temple Run	Arcade	Imangi Studios	4.170195
5	Subway Surfers	Arcade	Kiloo	4.496675
6	Hill Climb Racing	Racing	Fingetsoft	4.423215
7	Sonic Dash	Arcade	SEGA	4.49338
8	Super Stickman Golf 2	Sports	Noodlecake Studios Inc	4.41765

3. Python script:

```
merger.py x
1 import re
2
3 #all_predictions.csv is contains predicted output for all potential matches.
4 #In this script, we are checking for each potential match whether it is
5 #predicted as match or not. If it is matched then we are copying both tuples
6 #in matched_tuples.csv. We are merging both tuples according to merger rules
7 #and copying the merged tuple in merged_tuples.csv.
8 fq = open('all_predictions.csv', 'r')
9 fp = open('merged_tuples.csv', 'w')
10 ft = open('matched_tuples.csv', 'w')
11
12 #first line in merged_tuples.csv file for providing headings.
13 a = '"ID",'+'"name",'+'"category",'+'"developer",'+'"rating"'
14 fp.write(a)
15 fp.write('\n')
16
17 #first line in matched_tuples.csv file for providing headings.
18 #Each line in this file contains matched tuples from both google table and apple table.
19 a = '"G_ID",'+'"G_name",'+'"G_category",'+'"G_developer",'+'"G_rating",'+\
20     '"A_ID",'+'"A_name",'+'"A_category",'+'"A_developer",'+'"A_rating"'
21 ft.write(a)
22 ft.write('\n')
23
24 #initializing j here which will act as ID.
25 j = 0
26
27 for line in fq:
28     word = line.split()
29     #checking if the rows are predicted as matched or not.
30     if (word[len(word)-1] == "1"):
31         #writing ID in the first field
32         j = j + 1
33         idx = '">'+str(j)+'"'
34         fp.write(idx)
35
36         #extracting row number of each tuple which are being matched.
37         word_google = word[2]
38         word_apple = word[3]
39
40         #variables to hold attribute values of google table
41         ID_google = ""
42         name_google = ""
43         category_google = ""
44         developer_google = ""
45         rating_google = ""
46
47         #variables to hold attribute values of apple table
48         ID_apple = ""
49         name_apple = ""
50         category_apple = ""
51         developer_apple = ""
52         rating_apple = ""
53
54         with open('google.csv') as fr:
55             with open('apple.csv') as fs:
56                 #based on the extracted row number, reading the corresponding row
57                 #from google table and writing whole row in matched_tuples.csv file.
58                 content_google = fr.read().splitlines()
59                 row = content_google[int(word_google[1:-1])]
60                 ft.write(row)
61                 ft.write(",")
62
63                 #extracting different attribute values for merging.
64                 #extracting ID.
65                 start = row.find('"') + 1
66                 end = row.find('"', start)
67                 ID_google = row[start:end]
68
```

```

69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142

#extracting name.
row = row[(end+1):]
start = row.find('"') + 1
end = row.find('"', start)
name_google = row[start:end]

#extracting category.
row = row[(end+1):]
start = row.find('"') + 1
end = row.find('"', start)
category_google = row[start:end]

#extracting developer.
row = row[(end+1):]
start = row.find('"') + 1
end = row.find('"', start)
developer_google = row[start:end]

#extracting rating.
row = row[(end+1):]
start = row.find('"') + 1
end = row.find('"', start)
rating_google = row[start:end]

#based on the extracted row number, reading the corresponding row
#from apple table and writing whole row in matched_tuples.csv file.
content_apple = fs.read().splitlines()
row = content_apple[int(word_apple[1:-1])]
ft.write(row)
ft.write('\n')

#extracting different attribute values for merging.
#extracting ID.
start = row.find('"') + 1
end = row.find('"', start)
ID_apple = row[start:end]

#extracting name.
row = row[(end+1):]
start = row.find('"') + 1
end = row.find('"', start)
name_apple = row[start:end]

#extracting category.
row = row[(end+1):]
start = row.find('"') + 1
end = row.find('"', start)
category_apple = row[start:end]

#extracting developer.
row = row[(end+1):]
start = row.find('"') + 1
end = row.find('"', start)
developer_apple = row[start:end]

#extracting rating.
row = row[(end+1):]
start = row.find('"') + 1
end = row.find('"', start)
rating_apple = row[start:end]

#if both google table and apple table contain rating then taking
#average of both ratings. Otherwise, taking rating from google table.
if((rating_apple.find('Ratings') == -1) and (rating_apple.find('none') == -1)):
    rating_google = (float(rating_google) + float(rating_apple))/2

#combining all attributes as a single string and writing it in merged_tuples.csv.
to_write = ''+name_google+'',''+category_google+'',''+\
    ''+developer_google+'',''+str(rating_google)+''
fp.write(to_write)
fp.write('\n')

fp.close()
fq.close()
ft.close()

```