

Data Science [ITE4005]

**Programming Assignment #1**

: find association rules using the Apriori algorithm

2015004120

Software Engineering

Park Seonha

## 0. Environment

- OS : Windows
- Language : Python3 (version : 3.4.4)
- Executable file : apriori.py (./Programming\_Assignment\_1/project\_apriori/apriori.py)

## 1. Summary of Algorithm

This programming assignment's goal is "find association rules using the Apriori algorithm." So I approached this assignment in two parts, "implement of apriori algorithm" and "make association rules with all frequent itemsets."

To implement apriori algorithm, I referenced the pseudo-code of apriori algorithm on class slides. This algorithm takes following steps, and I will explain this steps with small example which the number of transactions is 4 and minimum support is 2.

Ⓐ make  $C_1$  (1 - candidate itemset) by scanning all transactions.

1	3	4		Scanning all transactions and count supports of each item_id.  $C_1 = [\{1\}: 2, \{2\}: 3, \{3\}: 3, \{4\}: 1, \{5\}: 3]$
2	3	5		
1	2	3	5	
2	5			
small example				

Ⓑ make  $L_1$  (1 - frequent itemset) by eliminate infrequent itemsets in  $C_1$

$C_1 = [\{1\}: 2, \{2\}: 3, \{3\}: 3, \{4\}: 1, \{5\}: 3]$ , minimum support = 2
Since minimum support is 2, {4} is infrequent itemset and others are frequent itemset. Add frequent itemsets in empty $L_1$ and the result is, $L_1 = [\{1\}: 2, \{2\}: 3, \{3\}: 3, \{5\}: 3]$

Ⓒ after k is 2, make  $C_k$  using  $L_{k-1}$  (do self-join and pruning infrequent itemsets)

$L_1 = [\{1\}, \{2\}, \{3\}, \{5\}]$
$C_2 = L_1 \bowtie L_1 = [\{1,2\}, \{1,3\}, \{1,5\}, \{2,3\}, \{2,5\}, \{3,5\}]$ All subsets of itemset in $C_2$ are in $L_1$ , I don't have to prune infrequent itemset. and get supports of each itemset in $C_2$ , the result of this stage is $C_2 = [\{1,2\}: 1, \{1,3\}: 2, \{1,5\}: 1, \{2,3\}: 2, \{2,5\}: 3, \{3,5\}: 2]$

Ⓓ make  $L_k$  pick up frequent itemsets in  $C_k$

$C_2 = [\{1,2\}: 1, \{1,3\}: 2, \{1,5\}: 1, \{2,3\}: 2, \{2,5\}: 3, \{3,5\}: 2]$ , minimum support = 2
{1,2} and {1,5} are not frequent set, so eliminate those itemsets. $L_2 = [\{1,3\}: 2, \{2,3\}: 2, \{2,5\}: 3, \{3,5\}: 2]$

Ⓔ repeat Ⓒ and Ⓓ while  $L_{k-1}$  is not empty set.

$L_2 = [\{1,3\}, \{2,3\}, \{2,5\}, \{3,5\}]$ $C_3 = L_2 \bowtie L_2 = [\{1,2,3\}, \{1,3,5\}, \{2,3,5\}]$ $\{1,2,3\}$ is infrequent set because its subset $\{1,2\}$ is infrequent. [Apriori pruning principle] Same reason, $\{1,3,5\}$ is infrequent set because of its subset $\{1,5\}$ . $C_3 = [\{2,3,5\} : 2]$ , minimum support = 2 $\{2,3,5\}$ is frequent set so $L_3$ is, $L_3 = [\{2,3,5\} : 2]$
$L_3 = [\{2,3,5\}]$ $C_4 = L_3 \bowtie L_3 = \emptyset$ , so $L_4$ is empty set, and algorithm end.

If  $L_{k-1}$  is empty set,  $C_k$  will be empty set and  $L_k$  will be also empty, so this algorithm ends in terminal condition [ $L_{k-1}$  is empty set].

The result of this apriori algorithm is set of all frequent itemsets and supports.

$$L = [\{1\}: 2, \{2\}: 3, \{3\}: 3, \{5\}: 3, \{1,3\}: 2, \{2,3\}: 2, \{2,5\}: 3, \{3,5\}: 2, \{2,3,5\}: 2]$$

Second part is making association rule with set of frequent itemsets  $L$ . In this process, I don't need 1-frequent itemset since the formula of generating association rule is  $s \rightarrow (t-s)$  with non-empty set  $t$  and its non-empty subset  $s$ . So, I can remove 1-frequent itemsets, and the result is below.

$$L' = [\{1,3\}: 2, \{2,3\}: 2, \{2,5\}: 3, \{3,5\}: 2, \{2,3,5\}: 2]$$

For each itemset in  $L'$ , make all non empty subsets and its difference set. Also,  $P(t-s|s)$  is  $P(t)/P(s)$ , and it is each association rule's confidence.

$\{1,3\}$	$\{1\} \rightarrow \{3\}$	0.5	1
subset : $\{1\}, \{3\}$	$\{3\} \rightarrow \{1\}$	0.5	0.67
$\{2,3\}$	$\{2\} \rightarrow \{3\}$	0.5	0.67
subset : $\{2\}, \{3\}$	$\{3\} \rightarrow \{2\}$	0.5	0.67
$\{2,5\}$	$\{2\} \rightarrow \{5\}$	0.75	1
subset : $\{2\}, \{5\}$	$\{5\} \rightarrow \{2\}$	0.75	1
$\{3,5\}$	$\{3\} \rightarrow \{5\}$	0.5	0.67
subset : $\{3\}, \{5\}$	$\{5\} \rightarrow \{3\}$	0.5	0.67
$\{2,3,5\}$	$\{2\} \rightarrow \{3,5\}$	0.5	0.67
subset : $\{2\}, \{3\}, \{5\}, \{2,3\}, \{2,5\}, \{3,5\}$	$\{3\} \rightarrow \{2,5\}$	0.5	0.67
	$\{5\} \rightarrow \{2,3\}$	0.5	0.67
	$\{2,3\} \rightarrow \{5\}$	0.5	1
	$\{2,5\} \rightarrow \{3\}$	0.5	0.67
	$\{3,5\} \rightarrow \{2\}$	0.5	1

## 2. Description of codes

### get command line arguments

```
5 import sys
6
7 #to change raw input to command line argument
8 minsup_param = int(sys.argv[1])
9 input_file = sys.argv[2]
10 output_file = sys.argv[3]
11
12 minsup = minsup_param / 100
13
```

To see code top to bottom, before start main process, program get command line arguments(minimum support, input file name, output file name) with **sys** module.

### implemented funtions

```
14 # change item_set between list and string {[item_id],[item_id],...[item_id]}
15 def list_to_str(item_set):
16     brace = "{"
17     for item_id in item_set:
18         brace = brace + str(item_id)
19         brace += ","
20     brace = brace[:-1]
21     brace += "}"
22     return brace
23
24 def str_to_list(str_list):
25     str_list = str_list[1:-1]
26     new_list = list(map(int, str_list.split(",")))
27     return new_list
28
```

This code is written in Python3 and because of its specialty, I have to define some functions to implement this assignment. Since Python can't make [key : list, value : int] dictionary, I made functions which can change list to string and list-formed string to list. Since I implement this functions, candidate itemset and frequent itemset are implemented as [key : string(list-formed), value : int], and also available set operations.

*list\_to\_str(item\_set)* gets list as parameter, and return list-formed string. It is same as assignment item set output format. For example if *list\_to\_str*([1, 3, 4]) executed, it returns string "{1,3,4}".

In reverse, function *str\_to\_list(str\_list)* gets string parameter like "{3,7,8}". After remove braces("{", "}") and commas, this function returns list [3, 7, 8].

```
30 # make subset of itemset
31 def jin_subset(item_set):
32     result_set = [[]]
33     for x in item_set:
34         result_set.extend([y + [x] for y in result_set])
35     # this is empty set
36     result_set.pop(0)
37     # this is not proper-subset
38     result_set.pop(len(result_set)-1)
39     for it in result_set:
40         it.sort()
41     return result_set
42
```

*jin\_subset(item\_set)* function gets list as parameter, and returns list of all subset of parameter. If *jin\_subset([1,2,3])* executed, 33 line and 34 line makes all subset of *item\_set* ( $\{\emptyset, [1], [2], [1,2], [3], [1,3], [2,3], [1,2,3]\}$ ), and after 36 line and 38 line, this function removes empty set and set which is exactly same with parameter set. Now, there are only remain proper subset, return this set of proper subsets after sorting each proper subsets. In this case, this function returns  $\{[1], [2], [1,2], [3], [1,3], [2,3]\}$ .

```

44 # return 1 if set1 and set2 is same set, return 0 other case
45 def set_comp(set1, set2):
46     result = 0
47     if len(set(set1).difference(set(set2))) == 0:
48         if len(set(set2).difference(set(set1))) == 0:
49             result = 1
50     return result

```

*set\_comp(set1,set2)* gets two sets as parameter, and checks these two sets are exactly same or not. If  $A-B$  is empty set and  $B-A$  is also empty set, they are same set because all elements are in  $A \cap B$ . So this function designed to check  $set1-set2$  and if it has zero length, in other words, if it is empty set, check if  $set2-set1$  is empty set or not. This function returns 1 if two sets are same, other case, returns 0.

```

52 #change round function because python round is not sa-sa-o-ip
53 def new_round(num):
54     new_num = num*1000
55     #if num is 13.255 -> new num is 13255
56     over_num = int(new_num/1000)
57     first_num = int((new_num%1000)/100)
58     second_num = int((new_num%100)/10)
59     third_num = int(new_num%10)
60     upper = 0
61     if third_num >= 5:
62         second_num = second_num + 1
63     if second_num > 9:
64         second_num = second_num - 10
65         first_num = first_num + 1
66     if first_num > 9:
67         first_num = first_num - 10
68         upper = 1
69
70     ret_num = int(over_num) + int(upper) + round(float(first_num/10),1) + round(float(second_num/100),2)
71     ret_num = round(ret_num,2)
72     return ret_num

```

*new\_round(num)* is the function gets a floating point number as parameter and return its two decimal place rounded number. Since python internal round function implied 'round to nearest even' way and human uses 'round to nearest integer' way, I defined new round off function. In 'round to nearest integer' way, we don't use after three decimal place, so if input is  $ab.cdefg_{(10)}$ , after 54 line, it changes  $abcde.fg_{(10)}$  and separated  $ab_{(10)}$ (over\_num),  $c_{(10)}$ (first\_num),  $d_{(10)}$ (second\_num),  $e_{(10)}$ (third\_num).

round off third\_num and regrouping with second\_num, first\_num, and merge them as  $ab_{(10)} + c/10_{(10)} + d/100_{(10)} = ab.cd_{(10)}$ . After merging, to prevent long tail numbers comes from floating point number's inaccuracy like 13.100000000000000135, I use round

function at two decimal place and it doesn't affect return value's accuracy because we only left number on two decimal place, so under two decimal place never raise up.

### Import input data and parsing

```
75 # open input data file and store in transaction list
76 with open(input_file) as f:
77     input_data = f.readlines()
78     input_data = [d.strip() for d in input_data]
79
```

after 75 line, main process of this program runs. Open input file and read file as [str, ...] format. Since all transaction strings finish "\n", I used strip() remove newline.

```
80 # parse each input lines to item_ids and get the number of total transactions
81 transactions = []
82 total_trans = 0
83 for input_line in input_data:
84     trans_list = list(map(int, input_line.split("\t")))
85     transactions.append(trans_list)
86     total_trans += 1
87
```

Since transaction stored like "1Wt2Wt3Wt5", parse this string to list of integer. *transactions* is the list of integer list. and while appending each transaction, count total number of transactions(*total\_trans*).

After this process, this program will have [[item\_id, item\_id, ...], ...] list and the number of transactions, and they are necessary to execute apriori algorithm.

### Start apriori algorithm

```
88
89 # Start apriori algorithm
90 cand = {}
91 freq = {}
92
```

*cand* is the list which will contain  $C_1 \dots C_k$  in dictionary type. And *freq* is the list which will contain  $L_1 \dots L_k$  in dictionary type. So their structure is list of dictionary whose key is string(list-formed) and value is integer.

```
93 # make C_1 (1-candidate itemset) by scanning transactions
94 # since list cannot be key of dictionary in python3, I changed sorted itemset to string
95 cand_1 = {}
96 for trans in transactions:
97     for item_id in trans:
98         itemset = []
99         itemset.append(item_id)
100         itemset_key = list_to_str(itemset)
101         if itemset_key not in cand_1:
102             cand_1[itemset_key] = 1
103         else:
104             cand_1[itemset_key] += 1
105 cand.append(cand_1)
106
```

First step of apriori algorithm is making  $C_1$  (1- candidate set) by using *transactions*.

scan all transactions in *transactions*, make list *itemset*, and change *itemset* list to string to use as the key of dictionary in Python. (In python, list cannot be the key of dictionary because it is not hashable type) After make *itemset\_key*, check if  $C_1$  already has *itemset\_key*. If it exists, count up its value, doesn't exist, add new key and value to  $C_1$  whose key is *itemset\_key*, value is 1.

After generate dictionary  $C_1$ , append  $C_1$  into list *cand*. (*cand*[0] =  $C_1$ )

```

107 # make L_1 (1-frequent itemset) by pruning C_1
108 freq_1 = {}
109 for key, value in cand_1.items():
110     if value / total_trans >= minsup:
111         freq_1[key] = value
112 freq.append(freq_1)

```

$L_1$  is also dictionary whose key is string and value is integer. To generate  $L_1$  using  $C_1$ , check if itemset's support overs minimum support or not.

After generate dictionary  $L_1$ , also append  $L_1$  in *freq* list.

```

115 k = 2

117 # for(k=2; L_k != empty set; k++)
118 while k <= total_trans and len(freq[k - 2]) != 0:
119     # make cand_k(C_k) from freq_km(L_{k-1}) joining itself and pruning
120     cand_k = {}
121     freq_km = freq[k - 2]
122
123     # self-joining step
124     joined_fkm = []
125     for k1, v1 in freq_km.items():
126         for k2, v2 in freq_km.items():
127             kj_s = set(str_to_list(k1)).union(set(str_to_list(k2)))
128             kj_l = list(kj_s)
129             kj_l.sort()
130             if len(kj_l) != k:
131                 continue
132             else:
133                 if list_to_str(kj_l) not in joined_fkm:
134                     joined_fkm.append(list_to_str(kj_l))
135             else:
136                 continue

```

After make  $C_1$  and  $L_1$ , we can get  $C_k$  with using  $L_{k-1}$ . ( $k \geq 2$ ) This loop will execute while  $L_{k-1}$  is not empty set.

*cand\_k* is  $C_k$ , *freq\_km* is  $L_{k-1}$ . And we need to self-join  $L_{k-1}$ . Since we are generating k-candidate itemset, union two set in  $L_{k-1}$ , and translate it to list and remain sets whose size is  $k$ . If it can be candidate itemset, append that list in *joined\_fkm* as string format.

In this result, all candidate itemsets (whether it is frequent or not) are in *joined\_fkm*.

```

138     # pruning step
139     # check its subset in freq_km, freq_kmm ...
140     # make each joined_fkm's non-empty subset
141     pruned_fkm = []
142     for j_set in joined_fkm:
143         subs_fkm = jin_subset(str_to_list(j_set))
144         prune = -1
145         for subset in subs_fkm:
146             prune = 1
147             for fset in freq:
148                 for key_fset, val_fset in fset.items():
149                     if set_comp(str_to_list(key_fset), subset) == 1:
150                         prune = 0
151             if prune == 1:
152                 break
153         if prune == 1:
154             continue
155         else:
156             pruned_fkm.append(j_set)
157     for p_set in pruned_fkm:
158         cand_k[p_set] = 0

160     # for each transactions
161     # increment the count of each candidate in C_k which are contained in t
162     for candidate, value in cand_k.items():
163         cand_list = str_to_list(candidate)
164         for transact in transactions:
165             if len(list(set(cand_list).difference(set(transact)))) == 0:
166                 cand_k[candidate] += 1

```

After make *joined\_fkm*, we need to prune infrequent itemset. According to Apriori pruning principle, if a subset of set *S* is infrequent, *S* must be infrequent set. Using this principle, I check sets in *joined\_fkm* should be pruned or not. Check all non-empty subsets of each set in *joined\_fkm*, and if any subset is not in *freq*, that set will be pruned.

After pruning, make  $C_k$  with *pruned\_fkm*. Since *cand\_k* is key-value dictionary and I don't know each of their support, all candidate itemset's support initialized 0 and get counts check *transactions*. If a transaction includes candidate itemset (this can be checked by *candidate* — *transaction* is empty set or not), that candidate counts up.

```

168     # make freq_k(L_k) in cand_k(C_k) with minsup
169     freq_k = {}
170     for candidate, value in cand_k.items():
171         if cand_k[candidate] / total_trans >= minsup:
172             freq_k[candidate] = value
173

```

Now,  $C_k$  is generated. And generating  $L_k$  is available. Check all itemsets in  $C_k$  if its support is same or over minimum support. If it is, add that itemset in *freq\_k*.

```

174     cand.append(cand_k)
175     freq.append(freq_k)
176     k += 1

```

Append  $C_k$  and  $L_k$  to *cand* and *freq*. increase *k*, return to top of loop, and generate  $C_{k+1}$  and  $L_{k+1}$  during  $L_k$  is not empty set. if 'while'(line 118) ends, apriori algorithm part is end.



## making association rules

```
179 freq.pop(len(freq)-1)
```

Since while terminology condition is  $L_{k-1}$  is empty set,  $\text{freq}[\text{size of freq} - 1]$  is empty itemset dictionary. So remove last item of *freq*.

```
190 decomp_freq_dic = {}
191 decomp_freq = []
192 for fset in freq:
193     for key, value in fset.items():
194         if len(str_to_list(key)) > 1:
195             decomp_freq.append(str_to_list(key))
196             decomp_freq_dic[key] = value
197
```

Structure of *freq* is [{string:int}, ..., ...] (list of the list of dictionary whose key is string and value is integer). It is too complicate so it needs to be decomposed.

I decomposed *freq* to *decomp\_freq* (string array), *decomp\_freq\_dic* (dictionary whose key is string and value is integer). for all dictionary list in *freq*, append its itemset to decomposed list and dictionary. And 1-frequent itemset has no non-empty proper subset, I removed itemsets in  $L_1$  to check the size of list is bigger than 1.

```
198 #make association rules with each subset and write in output file
199 for fset in decomp_freq:
200     set_size = len(fset)
201     all_subset = jin_subset(fset)
202     for subset in all_subset:
203         lhs = subset
204         rhs = list(set(fset).difference(set(lhs)))
205         rhs.sort()
206
207         s = decomp_freq_dic[list_to_str(fset)] / total_trans
208         c = (decomp_freq_dic[list_to_str(fset)] / total_trans) / (decomp_freq_dic[list_to_str(list(lhs))] / total_trans)
209         s *= 100
210         c *= 100
211         out_f.write(list_to_str(list(lhs)) + "\t" + list_to_str(list(rhs)) + "\t" + str(new_round(s)) + "\t" + str(new_round(c)) + "\n")
```

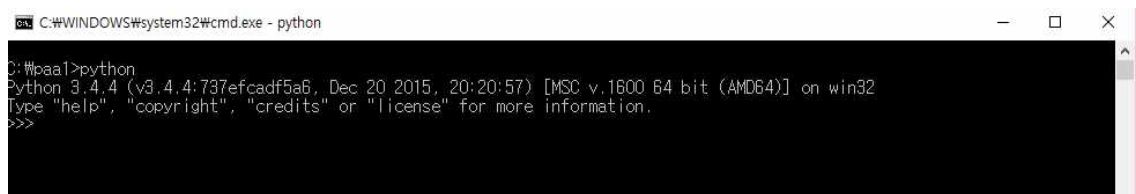
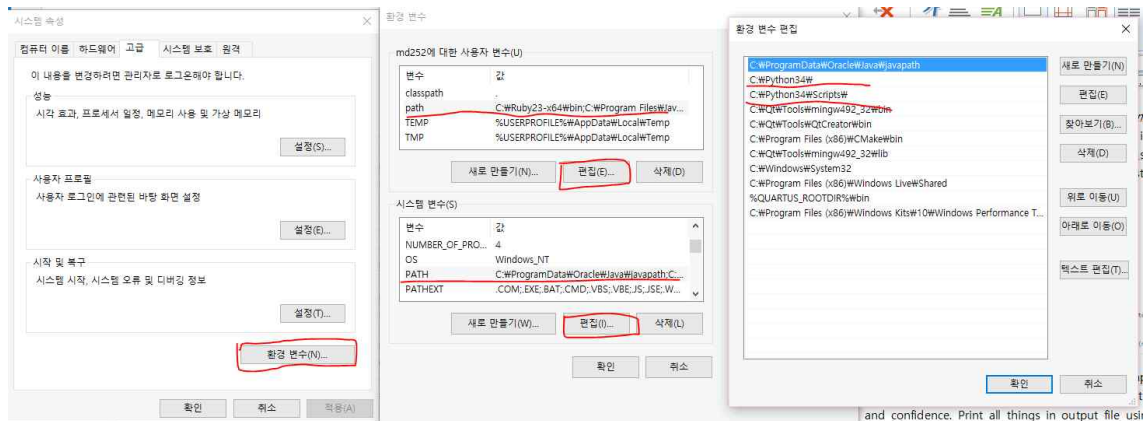
For frequent itemset in *decomp\_freq*, make all non-empty proper subset with function *jin\_subset()*, get difference set of each subset using set operation in python, calculate support(*s*) and confidence(*c*). Print all things in output file using *write()*. Since python *round()* function doesn't return value we expecting, I defined *new\_round()* function at top of code, and used here.

### 3. Instructions for compiling this code

This code is written in Python3 and tested in Python 3.4.4. To run this code, we need python 3.4.4 interpreter.

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

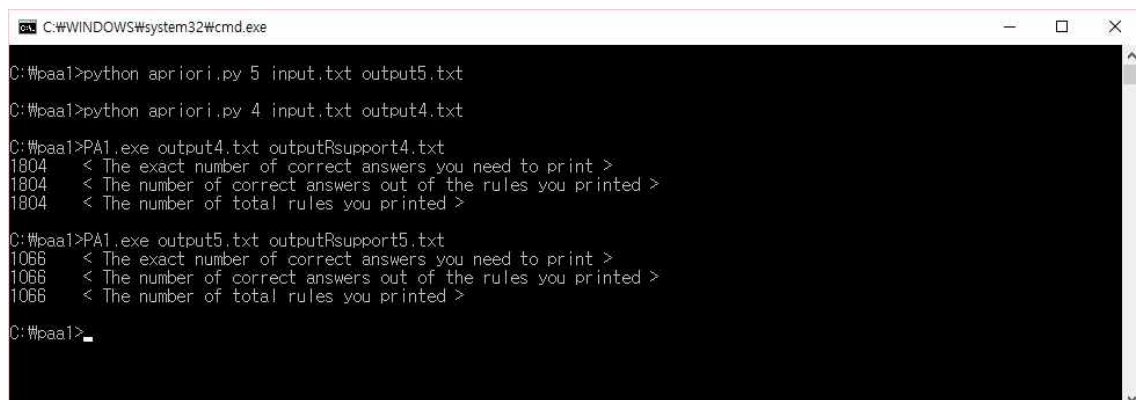
If python2 is already installed in PC, change path Python2x to Python34 in advanced system settings – environment variables.



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

In code directory(python file is in /project\_apriori), we can run this program on cmd "python apriori.py [minimum support] [input file name] [output file name]".

### 4. Other Specifications



In testing, some computers take quite long time to generate result because of their performance. (more than 40 seconds in minimum support 5%) It has additional condition in while loop ( $k$  is smaller than total\_trans), so it can't make infinite loop, just wait until it will finish.