

Apriori assignment report

2015004248
Jo Hyung Jung

1. Overview

‘Apriori’ is a candidate generation and test approach used in frequent pattern and association rule mining. By using ‘apriori pruning principle’, which denotes if there is any infrequent item set then its super set should not be generated, it can reduce many candidates and thus perform efficiently.

In this assignment, I got a chance that help me to deeply think about data mining method by implementing ‘apriori algorithm’.

2. Run-time environment

- OS: Ubuntu 18.04
- Language: Python 3.6
- Libraries:

Python sys module: for getting command-line arguments.

Python itertools module: for generating possible combinations in iterable object.

3. Summary of algorithm

‘Apriori’ algorithm consists of two stage repetition.

- i) Generate candidates by length [Self-Joining]
- ii) Prune candidates against data base [Pruning]

Self joining stage uses ‘apriori pruning principle’. By using this, it can reduce huge amount of candidates which is infrequent and thus perform faster.

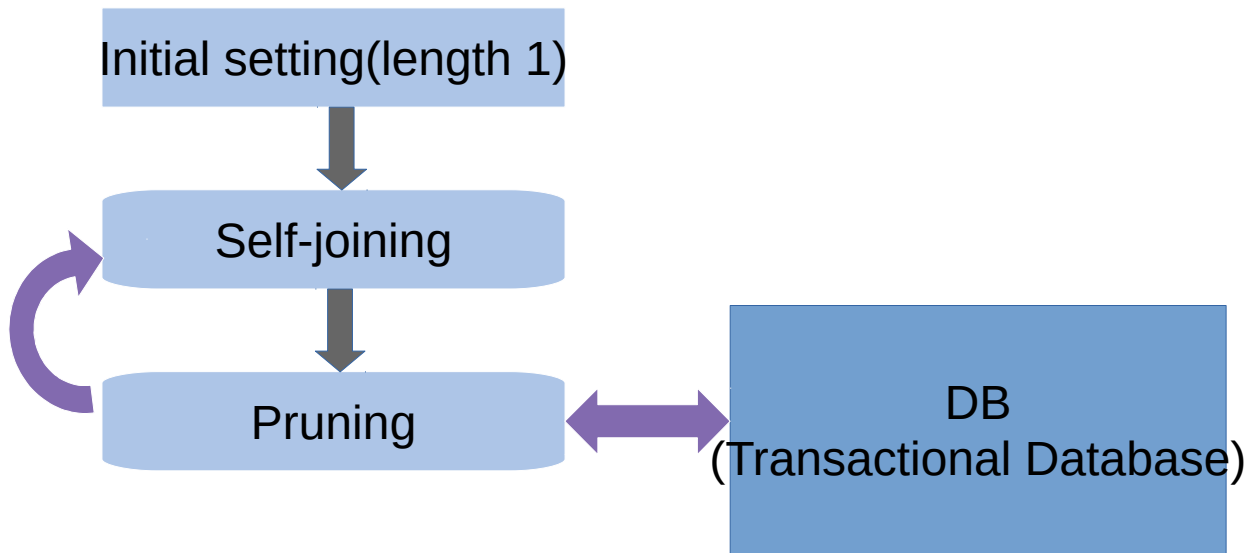
Pruning stage verifies candidates which generated in self-joining stage with minimum support. It needs scanning whole data base to calculate candidates’ support.

This two stage repetition continues until there are no other candidates. Verified item sets are saved with its support value.

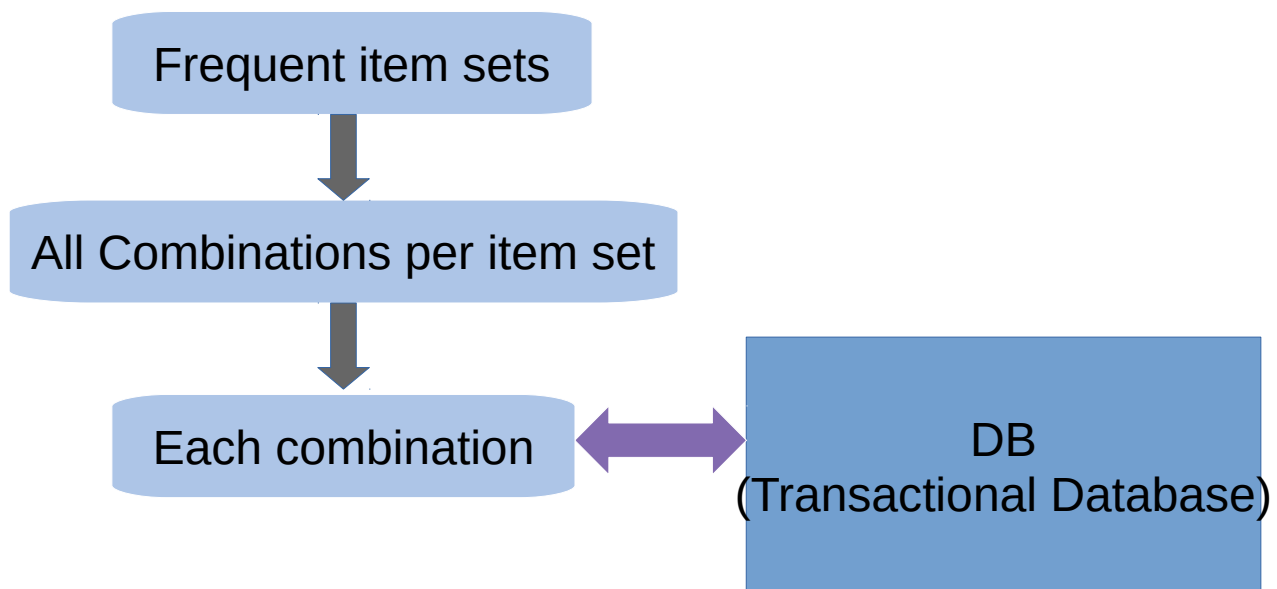
Also in this assignment, algorithm should generate proper association rules with its support and confidence. In this step, frequent item sets generate all possible combinations by each item.

Additional scanning in whole data base is required in each combination. Association rules are saved in proper format.

Below is the flow chart of this algorithm.



[Figure1: generating frequent item sets]



[Figure2: generating association rules]

i) Data Structures

- table: list of sorted sets

ex)

- **cand_list**: list of lists of sets

ex)

- freq_list: list of dictionaries(key: tuple, value: float)

ex)

$$[(1,): 50.0, (2,): 75.0, (3,): 75.0, (5,): 75.0], \{(1, 3): 50.0, (2, 3): 50.0, (2, 5): 75.0, (3, 5): 50.0\}, \{(2, 3, 5): 50.0\}]$$

ii) Functions

I use 3 functions to perform self-joining, pruning, and generating association rules.

- prune

```
4 '''
5     Feature: Prune candidates by support and length
6     Arguments: cand(type: list[set])
7     Return value type: dictionary(key: tuple, value: float)
8 '''
9 def prune(cand):
10     # Make counting table
11     cnt_table = {idx: 0 for idx, _ in enumerate(cand)}
12     for transaction in table:
13         for idx, itemset in enumerate(cand):
14             # Check whether each transaction is superset of each itemset
15             if transaction >= itemset:
16                 cnt_table[idx] += 1
17     # Filter by minimum support
18     ret = {tuple(cand[key]): val/len(table)*100 for key, val in cnt_table.items() if val>=min_sup}
19     return ret
20
21
```

Feature: Prune candidates by support and length

Arguments: cand → list of sets

Return value: ret → dictionary(key: tuple, value: float)

In this function, validate process is executed.

First, create counting table by using given candidates list. I use dictionary as this table, key for index of each candidate and value for its counts.

Second, scan the whole data base to count each candidate. Nested for loop is used.

Finally, filter the candidates with minimum support. I use python filter built-in function to compute easily. However, set can't be used as a key in dictionary, I convert its type set to tuple. Its support is stored as value.

- self_join

```
23 '''
24     Feature: Generate candidates by self-joining
25     Arguments: itemdict(type: dictionary(key: tuple, value: float)), length(type: integer)
26     Return value type: list[set]
27 '''
28 def self_join(itemdict,length):
29     # Convert type: tuple->set
30     itemsets = list(map(set,itemdict.keys()))
31     # Generate candidates by set union operator
32     ret = [itemsets[i]|itemsets[j] for i in range(len(itemsets)) for j in range(i+1,len(itemsets))]
33     # Filter by length
34     ret = list(filter(lambda x: len(x)==length, ret))
35     # Remove duplicates
36     ret = list(set(tuple(sorted(itemset)) for itemset in ret))
37     # Convert type: tuple->set
38     ret = list(map(set,ret))
39     return ret
```

Feature: Generate candidates by self-joining

Arguments: itemdict → dictionary(key: tuple, value: float) / length → integer

Return value: ret → list of sets

In this function, self-joining process is executed.

First, convert argument's type tuple to set. By doing this, set operator '|' (union) can be used to create candidates easily.

Second, generate all possible candidates by using set union operator.

Third, filter generated candidates with specific length. I use python built-in filter function to compute easily.

Finally, remove duplicates using some python tricks. As 'set' data structure doesn't allow duplicates, I use 'set' to remove duplicates. By applying 'set' to sorted tuple, all possible duplicates are removed easily.

- generate_association_rule

```
42 '''
43 Feature: Generate association rules
44 Arguments: freq_list(type: list[dictionary(key: tuple, value: float)])
45 Return value type: list[string]
46 '''
47 def generate_association_rule(freq_list):
48     ret = []
49     for itemdict in freq_list:
50         for itemset, val in itemdict.items():
51             # Exception: length is less or equal than 1
52             if len(itemset) <= 1:
53                 break
54             # Save support value
55             sup = val
56             # Convert type: tuple->set
57             itemset = set(itemset)
58             for length in range(1, len(itemset)):
59                 # Generate all combinations
60                 total_combinations = combinations(itemset, length)
61                 for combination in total_combinations:
62                     # Convert type: tuple->set
63                     left = set(combination)
64                     right = itemset - left
65                     left_cnt = 0
66                     for transaction in table:
67                         if transaction >= left:
68                             left_cnt += 1
69                     # Calculate confidence
70                     conf = sup * len(table) / left_cnt
71                     # Save result as string
72                     ret.append("{}\t{}\t{: .2f}\t{: .2f}\n".format(str(left), str(right), sup, conf))
73     return ret
```

Feature: Generate association rules

Arguments: freq_list → list of dictionaries(key: tuple, value: float)

Return value: ret → list of strings

In this function, association rules are generated.

First, make all possible combinations by using python itertools.combinations function. Each item set's support value is also stored.

Second, divide each combination into left side and right side. Notice that each side is converted to 'set' type for '-' operator(complement) and '>=' operator(is_superset).

Third, calculate confidence by scanning data base. Support and confidence are calculated as [# of (A&B) / # of total transaction] and [# of (A&b) / # of A] in association rule A → B so we can easily compute confidence by using 'support' and '# of A'.

Finally, each formatted string is stored in list. In this format, each side is enclosed in curly braces. Also, support and confidence is rounded in 2 decimal places.

iii) Code Flow

```
76 if __name__ == "__main__":
77     '''
78     Initialize arguments
79     [Type]
80     min_sup: float(ratio)->integer(count)
81     input_file: string
82     output_file: string
83     '''
84     # Save command line argument
85     min_sup, input_file, output_file = sys.argv[1:]
86     # Convert type: string->float
87     min_sup = float(min_sup)/100
88
89     '''
90     Read input file
91     & Create transaction table
92     [Type]
93     table: list[sorted sets]
94     '''
95     with open(input_file, 'r') as f:
96         table = [set(sorted(map(int, line.split())))) for line in f.readlines()]
97
98     # Convert: ratio->count
99     min_sup = len(table)*min_sup
100
101     '''
102     Generate candidates(Self-joining)
103     & Prune candidates
104     => Generate frequent item sets depending on its length
105     [Type]
106     itmes: set
107     cand_list: list[list[set]]
108     freq_list: list[dictionary(key:tuple, value=float(support))]
109     max_length: integer
110     '''
111     # Create candidates: length 1
112     items = {item_id for transaction in table for item_id in transaction}
113     cand_list = [{item_id} for item_id in items]
114     freq_list = [prune(cand_list[0])]
115     max_length = 1
```

1. Save command-line arguments, modify minimum support's type
 2. Create data base by reading input file, modify minimum support(Ratio to actual counts).
 3. Create candidates (length 1)
 4. Verify created candidates using 'prune' function (length 1)
- Now, we are ready to get into loop.

```
116     # Create candidates depending on its length
117     while True:
118         cand_list.append(self_join(freq_list[-1], max_length+1))
119         tmp = prune(cand_list[-1])
120         if not tmp:
121             break
122         freq_list.append(tmp)
123         max_length += 1
124
125     '''
126     Generate association rules
127     [Type]
128     association_list: list[string]
129     '''
130     association_list = generate_association_rule(freq_list)
131     # Write results in output file
132     with open(output_file, 'w') as f:
133         for line in association_list:
134             f.write(line)
```

5. Repeat creating candidates(using 'self_join' function) and verifying them(using 'prune' function) until there is no other candidates.
6. Generate association rules using 'generate_association_rule' function
7. Write the results in to output file

5. Instructions for execution

```
jsense@jsense:~/2020_ITE4005_2015004248/project_apriori$ python3 apriori.py 5 input.txt output.txt
jsense@jsense:~/2020_ITE4005_2015004248/project_apriori$ python3 apriori.py 50 sample_input.txt sample_output.txt
```

Type 'python3 "minimum support" "input file name" "output file name"' in terminal.

6. Test result

To verify code, I make simple sample file that can be easily figured out. This sample is in our lecture note. (Convert 'A','B','C'... in to '1','2','3'...) I verified this output file by calculating its support and confidence manually.

{1}	{3}	50.00	100.00
{3}	{1}	50.00	66.67
{2}	{3}	50.00	66.67
{3}	{2}	50.00	66.67
{2}	{5}	75.00	100.00
{5}	{2}	75.00	100.00
{3}	{5}	50.00	66.67
{5}	{3}	50.00	66.67
{2}	{3, 5}	50.00	66.67
{3}	{2, 5}	50.00	66.67
{5}	{2, 3}	50.00	66.67
{2, 3}	{5}	50.00	100.00
{2, 5}	{3}	50.00	66.67
{3, 5}	{2}	50.00	100.00

After successful sample test, I tested actual input file. Here is the result.

{8, 3}	{16, 17}	5.80	22.48
{8, 17}	{16, 3}	5.80	48.33
{16, 3}	{8, 17}	5.80	23.02
{16, 17}	{8, 3}	5.80	44.62
{17, 3}	{8, 16}	5.80	76.32
{8, 16, 3}	{17}	5.80	24.17
{8, 16, 17}	{3}	5.80	65.91
{8, 17, 3}	{16}	5.80	85.29
{16, 17, 3}	{8}	5.80	93.55
{8}	{16, 3, 13}	7.40	16.37
{16}	{8, 3, 13}	7.40	17.45
{3}	{8, 16, 13}	7.40	24.67
{13}	{8, 16, 3}	7.40	25.00
{8, 16}	{3, 13}	7.40	24.50
{8, 3}	{16, 13}	7.40	28.68
{8, 13}	{16, 3}	7.40	51.39
{16, 3}	{8, 13}	7.40	29.37
{16, 13}	{8, 3}	7.40	53.62
{3, 13}	{8, 16}	7.40	82.22

I got 1066 association rules in given input text file.

7. Reference

- [1] R. Agrawal and R. Srikant. Fast Algorithms for Mining Association Rules. VLDB, 1994.