

# Assignment 1: Apriori Algorithm

2019094511 김준표

## 1. Summary of algorithm

Apriori 알고리즘은 DB를 스캔한 후 item set의 길이를 1씩 늘려가며 후보 item set을 정하여 frequent item set을 찾는 알고리즘이다. 후보 item set을 정하는 과정에서 self-join과 pruning을 실시하며 입력 받은 minimum support를 기준으로 후보 선정 또는 탈락을 실시한다.

해당 알고리즘은 후보 아이템이 더 이상 없을 때까지 후보 생성 및 frequent item set을 찾는 일을 반복한다. K 단계에서의 self-join에서는 해당 단계에서 존재하는 item set을 조합하여 (K+1) 길이의 후보 아이템을 선정한다.

이후 pruning을 실시하여 새로 생성된 후보 item의 조합이 K 단계에서의 item set에 없는 경우를 제거한다. 예를 들면 (K=2) 단계에서 {A,B} {B,C} {B,E} {C,E}가 있을 때 (K=3)에서 {B,C,E}는 가능하지만 {A,B,C}는 {A,B}가 없으므로 생성될 수 없다. 이와 같은 경우를 제거해주는 작업이 pruning이다.

위의 작업을 모두 시행하여 새로운 item set을 생성했다면 해당 item set의 빈도수를 카운트하여 다음 단계로 넘겨줄 DB를 만들어주면 한 사이클이 마무리된다.

## 2. Code description

```
81 def main():
82     minimum_support = int(sys.argv[1])
83     input_file = sys.argv[2]
84     output_file = sys.argv[3]
85
86     transactions = []
87     input_items = {}
88     frequent_items = []
89
90     # DB의 transaction 정보들을 transactions table에 담는다
91     with open(input_file, "r") as f:
92         while True:
93             line = f.readline()
94             if not line:
95                 break
96             str = line.split("\n")
97             # 각 줄에 있는 아이템을 입력받아 transaction에 넣어준다
98             items = list(map(int, str[0].split('\t')))
99             transactions.append(items)
100            # transaction의 아이템 입력받기
101            for item in items:
102                item = tuple([item])
103                if item in input_items:
104                    input_items[item] += 1
105                else:
106                    input_items[item] = 1
107        f.close()
108        # 퍼센트로 입력받은 minimum support 값을 count로 변환
109        minimum_support_cnt = float(minimum_support / 100) * len(transactions)
110        # 빈번 아이템 필터링
111        frequent_items.append(filtering(input_items, minimum_support_cnt))
112
113        apriori(transactions, frequent_items, minimum_support_cnt)
114        association_rule(output_file, frequent_items, transactions)
115
```

main 함수에서는 input file로부터 받은 item 정보들을 바탕으로 초기 DB를 생성하고 Apriori 함수

를 호출한다.

```
35 # Apriori 알고리즘
36 def apriori(transactions, frequent_items, minimum_support_cnt):
37     cnt = 1
38     k = 1
39     while True:
40         # 후보 아이템 생성
41         # Step 1: self-join
42         # Step 2: pruning
43         candidate_items = self_join(frequent_items[k-1], cnt)
44         # 다음 후보 아이템이 없는 경우 종료
45         if len(candidate_items) == 0:
46             break
47         # 아이템과 빈도수를 저장하기 위한 딕셔너리
48         item_cnt_dict = {}
49         for item_set in transactions:
50             for item in candidate_items:
51                 tmp_item = tuple(item)
52                 if set(tmp_item) <= set(item_set):
53                     if tmp_item in item_cnt_dict:
54                         item_cnt_dict[tmp_item] += 1
55                     else:
56                         item_cnt_dict[tmp_item] = 1
57         frequent_items.append(filtering(item_cnt_dict, minimum_support_cnt))
58         cnt += 1
59         k += 1
60
```

Apriori 함수에서는 Apriori 알고리즘의 전체 틀을 바탕으로 동작하며 [frequent\_items]를 바탕으로 후보 item set인 [candidate\_items]를 추출하고, [candidate\_items]의 길이가 0이 될 때 (후보 아이템이 존재하지 않을 때)까지 일련의 작업을 수행한다. 후보 item set을 구한 뒤 전체 DB를 돌면서 유효한 item set과 빈도수를 저장하고 다음 단계로 넘겨줄 [frequent\_items]를 생성한다.

frequent\_items는 각 index마다 (index+1)의 길이를 가지는 item set과 해당 item set의 빈도수를 저장하는 dictionary를 모은 list로 구성되어 있다. 그래서 frequent\_items[0]에는 길이 1의 item과 빈도수를 저장하고, frequent\_items[1]에는 길이 2의 item과 빈도수를 저장하는 식으로 구성된다.

```
23 # k+1 번의 후보 생성
24 def self_join(items, cnt):
25     tmp_items = copy.deepcopy(items)
26     candidate = []
27     for set1 in tmp_items:
28         for set2 in tmp_items:
29             # 후보 선정을 원활히 하기 위해 union한 결과를 sort하여 후보 군에 추가
30             tmp = sorted(list(set(set1).union(set(set2))))
31             if len(tmp) == (cnt + 1) and tmp not in candidate:
32                 candidate.append(tmp)
33     return pruning(candidate, tmp_items, cnt)
```

Self\_join 함수에서는 K 단계의 item set을 입력 받아 해당 item set의 item들끼리 self-join을 수행한 후 (K+1)의 길이를 가지는 경우 후보군에 추가해준다. Self-join을 구현하기 위해 동일한 item set을 이중 for문으로 돌면서 각 set의 합집합인 union을 실시하고 검사하는 방식으로 후보를 선정했다.

```
10 # k 번의 아이템에 없는 요소를 갖는 아이템 삭제
11 def pruning(candidate, prev_list, prev_cnt):
12     for item in candidate:
13         sub_item = list(itertools.combinations(item, prev_cnt))
14         flag = True
15         for sub in sub_item:
16             if sub not in prev_list:
17                 flag = False
18                 break
19         if flag is False:
20             del candidate[candidate.index(item)]
21     return candidate
```

Pruning 함수는 self-join을 거쳐 생성된 후보 중에서 K 단계에서 없는 조합 요소를 가지는 item 을 삭제해준다. 이를 구현하기 위해 itertools.combination을 이용했으며 해당 모듈은 리스트의 요소를 받아서 parameter로 입력 받은 길이의 조합을 뽑아내는 함수이다.

예를 들어 test = ['A', 'B', 'C']의 리스트가 있고, c = itertools.combinations(test, 2)를 실시한다면, c는 다음의 결과를 갖는다. [('A', 'B'), ('A','C'), ('B','C')]

```
5 # 입력받은 아이템 중에서 minimum support 값을 넘는 요소만 반환
6 def filtering(item_db, min_sup):
7     frequent = dict(filter(lambda val: val[1] >= min_sup, item_db.items()))
8     return frequent
9
```

Filtering 함수는 minimum support 값을 넘는 요소만 반환해준다. 이를 구현하기 위해 lambda 함수를 이용해서 간단하게 {아이템: 빈도수}로 구성된 dictionary의 value를 기준으로 필터링을 실시했다.

```
64 # Association rule을 구하고 출력하는 함수
65 def association_rule(output_file, frequent_items, transactions):
66     t_num = len(transactions)
67     fo = open(output_file, "w")
68     output = ""
69     for i in range(len(frequent_items)): # 모든 행들을 순회하며
70         for item in frequent_items[i]: # 각 행에 대해서
71             for size in range(1, i+1): # 길이가 size인 조합 변환
72                 tmp_item = list(itertools.combinations(item, size))
73                 for tmp in tmp_item: # 생성된 조합들로 이루어진 집합의 모든 부분집합들을 비교
74                     result_item = list(item)
75                     result_item = tuple([x for x in result_item if x not in tmp])
76                     # X, Y 두 아이템이 같이 나오는 비율
77                     support = (frequent_items[i][item]/t_num) * 100
78                     # X가 나온 경우 X와 Y가 같이 나오는 비율
79                     item_set = format(tmp)
80                     associative_item_set = format(result_item)
81                     confidence = (frequent_items[i][item]/frequent_items[len(set(tmp))-1][tmp]) * 100
82                     output += f"{item_set}\t{associative_item_set}\t{str("%.2f"%round(support,2))+"\t"+str("%.2f"%round(confidence,2))+"\n"}
83     fo.write(output)
84     fo.close()
```

Association rule 함수에서는 입력 받은 frequent\_items 리스트에 저장된 모든 item 조합 요소와 associative item set을 출력하고, 이에 대한 support와 confidence를 구하여 출력한다.

```
61 def format(item_set):
62     return '{'+','.join(map(str, item_set))+ '}'
```

Association rule 함수 내에서 출력을 할 때, 출력 format을 맞추기 위해 format함수를 만들어줬다. 해당 함수를 통해 {[item],[item],[item]}과 같은 format을 지켜 출력하도록 했다.

### 3. Compiling method

실행환경:

OS: Windows 10

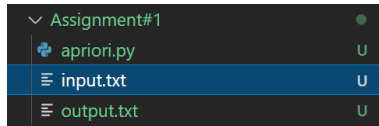
Python version: Python 3.8.3

실행방법:

```
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#1>
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#1> python .\apriori.py 5 input.txt output.txt
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#1>
```

python apriori.py 5 input.txt output.txt

위의 실행 명령어를 입력하면 아래와 같이 output.txt 파일이 생성되는 것을 확인할 수 있다.



#### 4. Other specification

```

Assignment#1 > input.txt
1 7 14
2 9
3 18 2 4 5 1
4 1 11 15 2 7 16 4 13
5 2 1 16
6 15 7 6 11 18 9 12 19 14
7 11 2 13 4
8 11 13
9 7 4 2 17 19 3 8 16 1
10 18 16 15 10 2 8 6 0 4 5

```

output.txt - Windows 메모장

파일(F)	편집(E)	서식(O)	보기(V)	도움말(H)
{7}	{14}	7.60	31.67	
{14}	{7}	7.60	29.69	
{2}	{18}	8.60	32.58	
{18}	{2}	8.60	31.16	
{4}	{18}	7.80	31.71	
{18}	{4}	7.80	28.26	
{5}	{18}	9.80	38.89	
{18}	{5}	9.80	35.51	
{1}	{18}	8.00	26.85	
{18}	{1}	8.00	28.99	
{2}	{4}	8.60	32.58	
{4}	{2}	8.60	34.96	
{2}	{5}	6.80	25.76	
{5}	{2}	6.80	26.98	
{1}	{2}	9.00	30.20	
{2}	{1}	9.00	34.09	
{4}	{5}	8.80	35.77	
{5}	{4}	8.80	34.92	
{1}	{4}	9.20	30.87	
{4}	{1}	9.20	37.40	
{1}	{5}	10.00	33.56	
{5}	{1}	10.00	39.68	
{2}	{7}	6.00	22.73	
{7}	{2}	6.00	25.00	
{4}	{7}	6.40	26.02	
{7}	{4}	6.40	26.67	

위와 같이 input.txt에 대해 output.txt가 출력되는 것을 확인할 수 있다.