

Data Mining Assignment2 - Frequent Itemsets

Georgios Leventis

November 2019

1 Introduction - Frequent Itemset Problem

One of the most frequent problems in the data mining field is the discovery and extraction of frequent items in sets of items(baskets). The problem can also be worded as finding itemsets that appear frequently in many of the same baskets. This is especially important for retail organizations that want to design and implement a customized marketing program[Fast Algorithms].

This problem is found in the "market-based" model of data. There is a large set of items $I = \{i_1, i_2, \dots, i_m\}$ and a large set of baskets $T = \{t_1, t_2, \dots, t_n\}$. Each basket is a subset of $T \subset I$.

The solution of the aforementioned problem aims to discover and create association rules, that are of the form $X \rightarrow Y$ (if X then Y, where X and Y are itemsets). Someone that bought $\{x, y, z\}$ usually buys $\{v, w\}$. Two more terms need to be explained:

- Support: Support for itemset I is the number of baskets containing all items in I. Given a threshold s, sets of items that appear in at least s baskets are frequent.
- Confidence: Confidence of the rule $I \rightarrow j$ is the ratio of the support for $I \cup \{j\}$ to the support for I, so: $\text{conf}(I \rightarrow j) = \frac{\text{support}(I \cup j)}{\text{support}(I)}$

2 Solving the Frequent Itemset Problem

There are some naive algorithms for solving the problem. However, they do not scale well if the dataset of the baskets is too large. This is where the Apriori algorithm comes in.

2.1 Apriori Algorithm

The Apriori Algorithm limits the memory demands by imposing a two-pass approach. During the 1st pass it reads the baskets and counts in the main memory the occurrences of each item. It requires $O(n)$ memory, with n = number of items. In the second pass it reads again the baskets and counts only

the pairs where both elements are frequent. It requires memory proportional to the square of frequent items; 2^m instead of 2^n .

3 Implementation

The implementation consists of the following functions:

- `ReadBaskets()`: This function reads the data and returns a list that contains all the baskets and the support, which is set to be 1% of the size of the list.

```

1  """
2
3  Returns baskets --> a list of all the baskets
4  """
5
6  def ReadBaskets():
7
8      infile = ('T10I40I00K.dat')
9      baskets = []
10
11      with open(infile, 'r') as file:
12          for transaction in file:
13              j = transaction.rstrip().split(" ")
14              basket = []
15              for item in j:
16                  item = int(item)
17                  basket.append(item)
18              baskets.append(basket)
19
20      support = int(0.01 * len(baskets))
21
22      return baskets, support
23
24 baskets, support = ReadBaskets()
25

```

Figure 1: `ReadBasket()` function

- `CreateItemWithSupport(baskets, support)`: Calculates the frequency that each items appears in the big list and returns a dictionary of singletons, with the item as the key and the number of appearances of this item in the list as the value.

```

26 """
27 Returns items_with_support --> dictionary with key = item and value = # of occurrences, i.e. (251 522)
28 """
29
30 def CreateItemWithSupport(baskets, support):
31     items_with_support = Counter(x for x in baskets)
32
33     for key in list(items_with_support):
34         if items_with_support[key] < support:
35             del items_with_support[key]
36     return items_with_support
37
38 items_with_support = CreateItemWithSupport(baskets, support)
39

```

Figure 2: `CreateItemsWithSupport()` function

- `CreateL1(items_with_support)`: Creates and returns a list that contains all the frequent singletons.

```

40 """
41 Creates the L1 set, set of the frequent singletons
42 """
43 def CreateL1(items_with_support):
44     L1 = set([])
45     for key in items_with_support:
46         L1.add(key)
47     return L1
48
49 L1 = CreateL1(items_with_support)

```

Figure 3: createL1() function

- `getLk(Lk_1, baskets, support, k)`: This function at first creates the C_k set of all the possible combinations, using the L_{k-1} set of frequent itemsets. Then it creates the L_k set of frequent items, after utilizing the `checkToAdd` function, to check if it can add a specific itemset to the set. Lastly, it deletes any itemset with support less than the threshold.

```

51 """
52 Creates the Ck -> list of all possible combinations from Lk_1 and uses its elements
53 to check if they are frequent.
54 """
55
56 def getLk(Lk_1, baskets, support, k):
57     Lk = {}
58
59     for basket in baskets:
60         bask = list(Lk_1.intersection(basket))
61         bask.sort()
62         Ck = list(combinations(bask, k))
63
64         for tuple in Ck:
65             if checkToAdd(tuple, Lk_1, k):
66                 if tuple in Lk:
67                     Lk[tuple] += 1
68                 else:
69                     Lk[tuple] = 1
70
71     for tuple in list(Lk):
72         if Lk[tuple] < support:
73             del Lk[tuple]
74
75     return Lk

```

Figure 4: getLk() function

- `checkToAdd(tuples, Lk_1, k)`: Checks if the itemsets in C_k actually exist in L_{k-1} . If they do it returns True, enabling the `getLk` function to add them to the set L_k .

```

51 """
52 Creates the Ck -> list of all possible combinations from Lk-1 and uses its elements
53 to check if they are frequent.
54 """
55
56 def getLk(Lk_1, baskets, support, k):
57     Lk = {}
58     for basket in baskets:
59         bask = list(Lk_1.intersection(basket))
60         bask.sort()
61         Ck = list(combinations(bask, k))
62
63         for tuple in Ck:
64             if checkToAdd(tuple, Lk_1, k):
65                 if tuple in Lk:
66                     Lk[tuple] += 1
67                 else:
68                     Lk[tuple] = 1
69
70     for tuple in list(Lk):
71         if Lk[tuple] < support:
72             del Lk[tuple]
73
74     return Lk

```

Figure 5: checkToAdd() function

- Apriori(L1, baskets, support): This function is initialized with $L_k = L_1$ and adds the frequent itemsets, regardless of the size, in a list. It continues to be executed, until the size of L_k is zero.

```

95
96 def Apriori(L1, baskets, support):
97     k = 1
98     Lset = L1
99     resultSet = [(1,)]
100     while(len(Lset) > 0):
101         resultSet.append(Lset)
102         print(resultSet)
103         Lset = getLk(Lset, baskets, support, k)
104         k += 1
105     return resultSet

```

Figure 6: apriori() function

4 Results

The results of the algorithm can vary, depending on the support threshold used. If the threshold is too large, then the program will produce only singletons and take little time to complete. As the threshold becomes smaller, then the algorithm will be able to produce itemsets with length more than one. The following figure depicts the results of the algorithm having as support 1% of the length of the dataset;1000: Lastly, Figure 8 shows the correlation between the

```

Please enter a number for the support: 1000
Number of itemsets with length 1: 375
Number of itemsets with length 2: 9
Number of itemsets with length 3: 1

```

Figure 7: Results with 1% support threshold

support threshold and the execution time of the algorithm. As seen, the larger the support, the smaller the execution time.

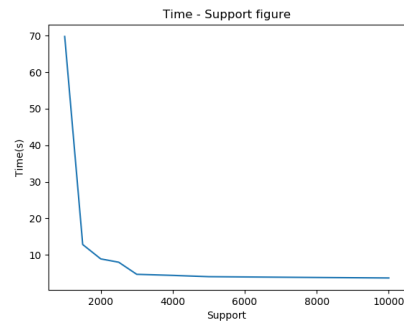


Figure 8: Support - Time correlation