

Q4 - Association Analysis

Consider the following market basket transactions shown in the Table below.

Transaction ID	Items Ordered
1	{Flour, Eggs, Bread}
2	{Soda, Coffee}
3	{Flour, Butter, Milk, Eggs}
4	{Bread, Eggs, Juice, Detergent}
5	{Bread, Milk, Eggs}
6	{Eggs, Bread}
7	{Detergent, Milk}
8	{Coffee, Soda, Juice}
9	{Butter, Juice, Bread}
10	{Milk, Bread, Detergent}

Import Libraries

```
1 from pprint import pprint
2 import itertools
3 import math
```

(a) How many items are in this data set? What is the maximum size of itemsets that can be extracted from this data set?

Transaction History

```

1 transactions = [(1, frozenset({'Flour', 'Eggs', 'Bread'})),
2                 (2, frozenset({'Soda', 'Coffee'})),
3                 (3, frozenset({'Flour', 'Butter', 'Milk', 'Eggs'})),
4                 (4, frozenset({'Bread', 'Eggs', 'Juice', 'Detergent'})),
5                 (5, frozenset({'Bread', 'Milk', 'Eggs'})),
6                 (6, frozenset({'Eggs', 'Bread'})),
7                 (7, frozenset({'Detergent', 'Milk'})),
8                 (8, frozenset({'Coffee', 'Soda', 'Juice'})),
9                 (9, frozenset({'Butter', 'Juice', 'Bread'})),
10                (10, frozenset({'Milk', 'Bread', 'Detergent'})),
11                ]

```

Count of Items

```

1 items = set()
2 for i, x in transactions:
3     items.update(x)
4 print('Total items in the data set - {}'.format(len(items)))
5 print('The items are - {}'.format(items))

```

```

1 Total items in the data set - 9
2 The items are - {'Juice', 'Flour', 'Detergent', 'Eggs', 'Bread', 'Milk',
   'Soda', 'Butter', 'Coffee'}

```

Itemsets with Support ≥ 1

```

1 itemsets = {}
2 for i in range(1, len(items)+1):
3     for x in itertools.combinations(items, i):
4         for _, y in transactions:
5             if set(x).issubset(y):
6                 if frozenset(x) not in itemsets:
7                     itemsets[frozenset(x)] = 0
8                 itemsets[frozenset(x)] += 1
9 print('Total itemsets with support  $\geq 1$  are - {}'.format(len(itemsets)))

```

```

1 Total itemsets with support  $\geq 1$  are - 44

```

```

1 print('Maximum Size of Itemset - {}'.format(max([len(x) for x in
   itemsets.keys()])))

```

There are 9 unique items, and 44 itemsets with support ≥ 1 , with maximum number of items in an itemset being 4.

(b) What is the maximum number of association rules that can be extracted from this data (including rules that have zero support)?

```
1 cnt=math.pow(3,len(items))-math.pow(2,len(items)+1)+1
2 print('The maximum number of association rules that can be extracted are
  {:.0f}'.format(cnt))
```

```
1 The maximum number of association rules that can be extracted are 18660
```

The maximum number of association rules that can be generated using d distinct items is given by the following formula - $R = \sum_{k=1}^{d-1} \left[\binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right] = 3^d - 2^{d+1} + 1$. If we plugin $d = 9$, $R = 18660$.

(c) What is the maximum number of 2-itemsets that can be derived from this data set (including those have zero support)?

```
1 print('The maximum number of 2-itemsets that can be derived from this data
  set are {:.0f}'.format(
2     math.factorial(len(items))/(math.factorial(len(items)-2)*2)))
```

```
1 The maximum number of 2-itemsets that can be derived from this data set are
  36
```

The maximum number of 2-itemsets that can be derived using n distinct items is given by the following combination - $count = \binom{n}{2} = \frac{n!}{n!2!} = \frac{n(n-1)}{2}$. For 9 items, we get 36 2-itemsets.

(d) Find an itemset (of size 2 or larger) that has the largest support.

```
1 max_cnt=0
2 max_set=None
3 for x,cnt in itemsets.items():
4     if len(x)>=2 and cnt>max_cnt:
5         max_cnt=cnt
6         max_set=x
7 print('The itemset of size 2 or larger with maximum support is -> {' + ',
      '.join(max_set) + '})')
```

```
1 The itemset of size 2 or larger with maximum support is -> {Bread, Eggs}
```

We can calculate this by iterating over all the itemsets, and checking itemsets which have support more than or equal to 2. We can then maintain a MAX variable to keep track of the set which has the maximum support. Using this approach we get the following itemset - > {Bread, Eggs} and its support is 4 or $\frac{4}{10} = 0.4$.

(e) Given minconf = 0.5, find two pairs of items, a and b, such that the rules {a} -> {b} and {b} -> {a} have the same confidence, and their confidence is greater than or equal to the minconf threshold.

```
1 s = set(itemsets.keys())
2 for x, y in itertools.combinations(s, 2):
3     union = x.union(y)
4     if union in itemsets and itemsets[x]==itemsets[y] and
       itemsets[union]/itemsets[x]>=0.5:
5         print('Confidence -> {:.2f}\t'.format(itemsets[union]/itemsets[x]),
6               '[]'.format(', '.join(x)),
7               '[]'.format(', '.join(y)))
```

```
1 Confidence -> 1.00    [Butter, Eggs] [Milk, Butter, Eggs]
2 Confidence -> 1.00    [Butter, Eggs] [Flour, Milk, Eggs]
3 Confidence -> 1.00    [Butter, Eggs] [Flour, Milk, Butter, Eggs]
4 Confidence -> 1.00    [Butter, Eggs] [Flour, Milk]
5 Confidence -> 1.00    [Butter, Eggs] [Flour, Butter, Eggs]
6 Confidence -> 1.00    [Butter, Eggs] [Flour, Butter]
```

7	Confidence -> 1.00	[Butter, Eggs] [Flour, Milk, Butter]
8	Confidence -> 1.00	[Butter, Eggs] [Milk, Butter]
9	Confidence -> 1.00	[Milk, Butter, Eggs] [Flour, Milk, Eggs]
10	Confidence -> 1.00	[Milk, Butter, Eggs] [Flour, Milk, Butter, Eggs]
11	Confidence -> 1.00	[Milk, Butter, Eggs] [Flour, Milk]
12	Confidence -> 1.00	[Milk, Butter, Eggs] [Flour, Butter, Eggs]
13	Confidence -> 1.00	[Milk, Butter, Eggs] [Flour, Butter]
14	Confidence -> 1.00	[Milk, Butter, Eggs] [Flour, Milk, Butter]
15	Confidence -> 1.00	[Milk, Butter, Eggs] [Milk, Butter]
16	Confidence -> 1.00	[Eggs, Detergent] [Bread, Juice, Eggs]
17	Confidence -> 1.00	[Eggs, Detergent] [Bread, Juice, Detergent]
18	Confidence -> 1.00	[Eggs, Detergent] [Juice, Eggs]
19	Confidence -> 1.00	[Eggs, Detergent] [Eggs, Juice, Detergent]
20	Confidence -> 1.00	[Eggs, Detergent] [Eggs, Bread, Detergent]
21	Confidence -> 1.00	[Eggs, Detergent] [Juice, Detergent]
22	Confidence -> 1.00	[Eggs, Detergent] [Eggs, Juice, Bread, Detergent]
23	Confidence -> 0.50	[Milk, Eggs] [Flour]
24	Confidence -> 0.50	[Milk, Eggs] [Bread, Milk]
25	Confidence -> 0.50	[Milk, Eggs] [Flour, Eggs]
26	Confidence -> 0.50	[Milk, Eggs] [Butter]
27	Confidence -> 1.00	[Flour] [Flour, Eggs]
28	Confidence -> 0.50	[Flour] [Butter]
29	Confidence -> 1.00	[Bread, Juice, Eggs] [Bread, Juice, Detergent]
30	Confidence -> 1.00	[Bread, Juice, Eggs] [Juice, Eggs]
31	Confidence -> 1.00	[Bread, Juice, Eggs] [Eggs, Juice, Detergent]
32	Confidence -> 1.00	[Bread, Juice, Eggs] [Eggs, Bread, Detergent]
33	Confidence -> 1.00	[Bread, Juice, Eggs] [Juice, Detergent]
34	Confidence -> 1.00	[Bread, Juice, Eggs] [Eggs, Juice, Bread, Detergent]
35	Confidence -> 1.00	[Juice, Coffee] [Juice, Soda, Coffee]
36	Confidence -> 1.00	[Juice, Coffee] [Juice, Soda]
37	Confidence -> 1.00	[Bread, Juice, Detergent] [Juice, Eggs]
38	Confidence -> 1.00	[Bread, Juice, Detergent] [Eggs, Juice, Detergent]
39	Confidence -> 1.00	[Bread, Juice, Detergent] [Eggs, Bread, Detergent]
40	Confidence -> 1.00	[Bread, Juice, Detergent] [Juice, Detergent]
41	Confidence -> 1.00	[Bread, Juice, Detergent] [Eggs, Juice, Bread, Detergent]
42	Confidence -> 1.00	[Juice, Eggs] [Eggs, Juice, Detergent]
43	Confidence -> 1.00	[Juice, Eggs] [Eggs, Bread, Detergent]
44	Confidence -> 1.00	[Juice, Eggs] [Juice, Detergent]
45	Confidence -> 1.00	[Juice, Eggs] [Eggs, Juice, Bread, Detergent]
46	Confidence -> 1.00	[Eggs, Juice, Detergent] [Eggs, Bread, Detergent]
47	Confidence -> 1.00	[Eggs, Juice, Detergent] [Juice, Detergent]
48	Confidence -> 1.00	[Eggs, Juice, Detergent] [Eggs, Juice, Bread, Detergent]
49	Confidence -> 1.00	[Juice, Soda, Coffee] [Juice, Soda]
50	Confidence -> 1.00	[Flour, Bread, Eggs] [Flour, Bread]

51	Confidence -> 1.00	[Eggs, Bread, Detergent] [Juice, Detergent]
52	Confidence -> 1.00	[Eggs, Bread, Detergent] [Eggs, Juice, Bread, Detergent]
53	Confidence -> 1.00	[Flour, Milk, Eggs] [Flour, Milk, Butter, Eggs]
54	Confidence -> 1.00	[Flour, Milk, Eggs] [Flour, Milk]
55	Confidence -> 1.00	[Flour, Milk, Eggs] [Flour, Butter, Eggs]
56	Confidence -> 1.00	[Flour, Milk, Eggs] [Flour, Butter]
57	Confidence -> 1.00	[Flour, Milk, Eggs] [Flour, Milk, Butter]
58	Confidence -> 1.00	[Flour, Milk, Eggs] [Milk, Butter]
59	Confidence -> 1.00	[Soda, Coffee] [Soda]
60	Confidence -> 1.00	[Soda, Coffee] [Coffee]
61	Confidence -> 1.00	[Flour, Milk, Butter, Eggs] [Flour, Milk]
62	Confidence -> 1.00	[Flour, Milk, Butter, Eggs] [Flour, Butter, Eggs]
63	Confidence -> 1.00	[Flour, Milk, Butter, Eggs] [Flour, Butter]
64	Confidence -> 1.00	[Flour, Milk, Butter, Eggs] [Flour, Milk, Butter]
65	Confidence -> 1.00	[Flour, Milk, Butter, Eggs] [Milk, Butter]
66	Confidence -> 1.00	[Soda] [Coffee]
67	Confidence -> 1.00	[Bread, Juice, Butter] [Juice, Butter]
68	Confidence -> 1.00	[Bread, Juice, Butter] [Bread, Butter]
69	Confidence -> 0.50	[Bread, Juice] [Bread, Detergent]
70	Confidence -> 0.50	[Bread, Juice] [Butter]
71	Confidence -> 0.50	[Milk, Detergent] [Bread, Detergent]
72	Confidence -> 0.50	[Milk, Detergent] [Bread, Milk]
73	Confidence -> 0.50	[Bread, Detergent] [Bread, Milk]
74	Confidence -> 1.00	[Juice, Detergent] [Eggs, Juice, Bread, Detergent]
75	Confidence -> 0.50	[Flour, Eggs] [Butter]
76	Confidence -> 1.00	[Flour, Milk] [Flour, Butter, Eggs]
77	Confidence -> 1.00	[Flour, Milk] [Flour, Butter]
78	Confidence -> 1.00	[Flour, Milk] [Flour, Milk, Butter]
79	Confidence -> 1.00	[Flour, Milk] [Milk, Butter]
80	Confidence -> 1.00	[Juice, Butter] [Bread, Butter]
81	Confidence -> 1.00	[Flour, Butter, Eggs] [Flour, Butter]
82	Confidence -> 1.00	[Flour, Butter, Eggs] [Flour, Milk, Butter]
83	Confidence -> 1.00	[Flour, Butter, Eggs] [Milk, Butter]
84	Confidence -> 1.00	[Flour, Butter] [Flour, Milk, Butter]
85	Confidence -> 1.00	[Flour, Butter] [Milk, Butter]
86	Confidence -> 1.00	[Flour, Milk, Butter] [Milk, Butter]

The formula for confidence is given by $confidence(X \rightarrow Y) = \frac{support(X \cup Y)}{support(X)}$. Given 2 itemsets, x and y, the confidence(x->y) and confidence(y->x) is equal if and only if the their individual supports are equal.

We listed out above all the possible pair of itemsets have a confidence ≥ 0.5 and have the same confidence.

If we assume a and b are items, then the question states to find 2 such pairs of 1-itemsets. From the above list these are [{Flour} {Butter}] and [{Soda}{Coffee}].

If we assume a and b are 2-itemsets, then the questions asks us to find a pair of 2-itemsets that satisfy the above condition. From the above list, this is [{Butter, Eggs}, {Flour, Milk}]