

## Lab5: Frede Emnetu (100704524)

```
In [1]: 1 import numpy as np
        2 import matplotlib.pyplot as plt
        3 import time
        4 import pandas as pd
        5 from itertools import combinations
```

### useful functions

```
In [2]: 1 def createCk(Lk, k):
        2     cand_list = []
        3     len_Lk = len(Lk)
        4
        5     for i in range(len_Lk):
        6         for j in range(i + 1, len_Lk):
        7             L1 = list(Lk[i])[k - 2]
        8             L2 = list(Lk[j])[k - 2]
        9             L1.sort()
       10             L2.sort()
       11             if L1 == L2:
       12                 cand_list.append(Lk[i] | Lk[j])
       13     return cand_list
```

```
In [3]: 1 # hash functions
        2 def hash1(i, j, numBuckets):
        3     return (i*j) % numBuckets
        4
        5
        6 def hash2(i, j, numBuckets):
        7     return (i+j) % numBuckets
```

```
In [4]: 1 def output(pairs):
        2     for x in pairs:
        3         list_ = list(x)
        4         print('%d --> %d' % (list_[0], list_[1]))
```

```

In [5]: 1 def apyori(itemsets, min_support, numBuckets):
2         start = time.time()
3         C1 = []
4         for basket in itemsets:
5             for item in basket:
6                 if not [item] in C1:
7                     C1.append([item])
8         C1 = [set(x) for x in C1]
9
10        count = {}
11        freq_items = []
12        L1 = []
13        countBuckets1 = [0] * numBuckets
14        countBuckets2 = [0] * numBuckets
15
16        for basket in itemsets:
17            for item in C1:
18                if item.issubset(basket):
19                    candidate = frozenset(item)
20                    if candidate not in count:
21                        count[candidate] = 1
22                    else:
23                        count[candidate] += 1
24        pairs = list(combinations(basket, 2))
25        for pair in pairs:
26            listPair = list(pair)
27            numbers = [int(x) for x in listPair]
28            countBuckets1[hash1(numbers[0], numbers[1], numBuckets)] += 1
29
30        for i in range(0, len(countBuckets1) - 1):
31            if countBuckets1[i] / len(itemsets) >= min_support:
32                countBuckets1[i] = 1
33            else:
34                countBuckets1[i] = 0
35
36        for key in count:
37            support = count[key] / len(itemsets)
38            if support >= min_support:
39                freq_items.insert(0, key)
40                freq_items.insert(1, support)
41                L1.insert(0, key)
42
43        # print("Frequent Items: ", freq_items)
44
45        freq_pairs = []
46        L2 = []
47        C2 = []
48        final = []
49        count = {}
50        test = [list(x) for x in L1]
51
52        for i in test:
53            final.append(int(i[0]))
54        # construct second hashmap
55        for basket in itemsets:
56            pairs = list(combinations(basket, 2))

```

```

57     for pair in pairs:
58         listPair = list(pair)
59         numbers = [int(x) for x in listPair]
60         # if the pair hashes to frequent bucket in first hashmap and bot
61         # hash the pair to the second hashmap
62         if countBuckets1[hash1(numbers[0], numbers[1], numBuckets)] == 1
63             countBuckets2[hash2(numbers[0], numbers[1], numBuckets)] += 1
64
65     for i in range(0, len(countBuckets2) - 1):
66         if countBuckets2[i] / len(itemsets) >= min_support:
67             countBuckets2[i] = 1
68         else:
69             countBuckets2[i] = 0
70
71     for basket in itemsets:
72         pairs = list(combinations(basket, 2))
73         for pair in pairs:
74             listPair = list(pair)
75             numbers = [int(x) for x in listPair]
76             # if pair hashes to frequent buckets in both hashmaps and values
77             if countBuckets1[hash1(numbers[0], numbers[1], numBuckets)] == 1
78                 C2.append(pair)
79
80     for basket in itemsets:
81         for item in C2:
82             item = frozenset(item)
83             if item.issubset(basket):
84                 candidate = frozenset(item)
85                 if candidate not in count:
86                     count[candidate] = 1
87                 else:
88                     count[candidate] += 1
89
90     for key in count:
91         support = count[key] / len(itemsets)
92         if support >= min_support:
93             freq_pairs.insert(0, key)
94             freq_pairs.insert(1, support)
95             L2.insert(0, key)
96     end = time.time()
97     return [end - start, L2]
98
99     # print("Frequent Pairs: ", freq_pairs)

```

## Retail DataSet

```

In [6]: 1 retail = pd.read_csv("http://fimi.uantwerpen.be/data/retail.dat", delimiter="
        2 itemsets = retail.values.tolist()

```

## Remove Nan Values

```
In [7]: 1 retail_NN = []
        2 for x in itemsets:
        3     retail_NN.append([i for i in x if str(i) != 'nan'])
```

## Create Sections

```
In [8]: 1 sections = [2,5,10,15, 20]
        2 associationR = []
        3 for i, x in enumerate(sections):
        4     new_list = retail_NN[0:int(len(retail_NN)*(sections[i]/100))]
        5     associationR.append(new_list)
```

## Use implemented function

```
In [9]: 1 times = []
        2 pairslen = []
        3 actualpairs = []
        4 for x in range(5):
        5     print(x)
        6     results = apyori(associationR[x],0.01,2000)
        7     times.append(results[0])
        8     pairslen.append(len(results[1]))
        9     actualpairs.append(results[1])
```

```
0
1
2
3
4
```

## Pairs

```
In [10]: 1 for x in range(len(sections)):
          2     print()
          3     print()
          4     print('with %g percent of the dataset I got %d pairs.' % (sections[x], p
          5     print()
          6     print(output(actualpairs[x]))
```

with 2 percent of the dataset I got 84 pairs.

```
41 --> 1327
38 --> 1327
32 --> 65
32 --> 271
749 --> 39
340 --> 39
381 --> 38
48 --> 301
41 --> 101
48 --> 1327
1327 --> 39
41 --> 271
41 --> 338
32 --> 110
589 --> 39
48 --> 740
740 --> 39
664 --> 39
38 --> 271
48 --> 110
604 --> 39
41 --> 255
48 --> 264
48 --> 310
41 --> 225
41 --> 170
48 --> 475
475 --> 39
48 --> 438
438 --> 39
89 --> 41
65 --> 41
48 --> 371
48 --> 60
371 --> 39
60 --> 39
371 --> 38
310 --> 39
38 --> 286
105 --> 38
48 --> 237
271 --> 39
270 --> 39
264 --> 39
225 --> 39
```

```
48 --> 255
255 --> 39
48 --> 65
237 --> 39
65 --> 39
179 --> 39
32 --> 179
32 --> 38
32 --> 48
32 --> 39
48 --> 170
170 --> 39
170 --> 38
48 --> 147
147 --> 39
41 --> 110
110 --> 39
110 --> 38
89 --> 38
48 --> 101
48 --> 95
48 --> 89
101 --> 39
89 --> 39
48 --> 79
48 --> 41
48 --> 36
79 --> 39
32 --> 41
48 --> 49
48 --> 39
48 --> 38
41 --> 39
41 --> 38
38 --> 39
37 --> 38
41 --> 36
36 --> 39
36 --> 38
None
```

with 5 percent of the dataset I got 126 pairs.

```
225 --> 60
89 --> 60
522 --> 1715
2990 --> 3735
1715 --> 604
604 --> 3966
60 --> 189
48 --> 3735
48 --> 2990
3735 --> 39
48 --> 1121
1715 --> 39
2238 --> 39
```

```
32 --> 286
48 --> 2238
2238 --> 38
32 --> 101
352 --> 1859
9 --> 310
41 --> 1327
32 --> 65
32 --> 271
48 --> 1715
1198 --> 39
48 --> 749
749 --> 39
32 --> 749
48 --> 301
41 --> 604
41 --> 310
41 --> 301
32 --> 310
41 --> 101
48 --> 1327
1327 --> 39
48 --> 271
48 --> 677
32 --> 110
208 --> 41
1121 --> 39
352 --> 39
41 --> 117
348 --> 38
522 --> 39
48 --> 824
824 --> 39
48 --> 740
41 --> 740
740 --> 39
677 --> 39
48 --> 201
38 --> 271
48 --> 604
48 --> 110
604 --> 39
48 --> 592
592 --> 39
48 --> 522
48 --> 310
48 --> 270
48 --> 225
41 --> 225
41 --> 170
189 --> 38
48 --> 475
475 --> 39
48 --> 438
438 --> 39
89 --> 41
65 --> 41
```

```
48 --> 60
60 --> 39
48 --> 352
48 --> 338
48 --> 161
310 --> 39
38 --> 286
48 --> 237
271 --> 39
270 --> 39
225 --> 39
48 --> 255
255 --> 39
48 --> 65
237 --> 39
36 --> 237
65 --> 39
48 --> 189
189 --> 39
48 --> 179
179 --> 39
32 --> 179
32 --> 38
32 --> 48
32 --> 39
48 --> 170
170 --> 39
170 --> 38
48 --> 147
147 --> 39
41 --> 110
110 --> 39
110 --> 38
89 --> 38
48 --> 101
48 --> 94
48 --> 89
101 --> 39
94 --> 39
89 --> 39
48 --> 79
48 --> 41
48 --> 36
79 --> 39
32 --> 41
48 --> 49
49 --> 38
48 --> 39
48 --> 38
41 --> 39
41 --> 38
38 --> 39
37 --> 38
41 --> 36
36 --> 39
36 --> 38
None
```



with 10 percent of the dataset I got 124 pairs.

```
225 --> 60
89 --> 260
89 --> 60
1715 --> 604
604 --> 3966
60 --> 189
476 --> 38
48 --> 1121
1715 --> 39
2238 --> 39
32 --> 286
2238 --> 38
855 --> 39
32 --> 101
352 --> 1859
9 --> 310
41 --> 1327
48 --> 1859
1859 --> 39
32 --> 271
48 --> 1715
1198 --> 39
48 --> 749
749 --> 39
65 --> 9
48 --> 301
41 --> 604
41 --> 310
41 --> 301
32 --> 310
41 --> 101
48 --> 1327
1327 --> 39
48 --> 271
65 --> 201
48 --> 1146
1146 --> 39
32 --> 110
208 --> 41
1121 --> 39
352 --> 41
352 --> 39
249 --> 237
41 --> 117
32 --> 117
522 --> 39
48 --> 9
48 --> 846
824 --> 39
48 --> 740
19 --> 38
533 --> 39
48 --> 201
```

```
38 --> 271
48 --> 604
48 --> 110
41 --> 475
604 --> 39
48 --> 533
48 --> 522
48 --> 310
48 --> 270
48 --> 225
41 --> 225
41 --> 170
189 --> 38
41 --> 237
48 --> 475
475 --> 39
48 --> 438
438 --> 39
89 --> 41
65 --> 41
48 --> 60
60 --> 39
48 --> 352
48 --> 161
310 --> 39
38 --> 286
48 --> 237
271 --> 39
270 --> 39
225 --> 39
48 --> 255
255 --> 39
48 --> 65
237 --> 39
36 --> 237
65 --> 39
48 --> 179
179 --> 39
32 --> 179
32 --> 38
32 --> 48
32 --> 39
48 --> 170
170 --> 39
170 --> 38
48 --> 147
147 --> 39
41 --> 110
110 --> 39
110 --> 38
89 --> 38
48 --> 101
48 --> 94
48 --> 89
101 --> 39
94 --> 39
89 --> 39
```

```
48 --> 79
48 --> 41
48 --> 36
79 --> 39
32 --> 41
48 --> 39
48 --> 38
41 --> 39
41 --> 38
38 --> 39
37 --> 38
41 --> 36
36 --> 39
36 --> 38
None
```

with 15 percent of the dataset I got 123 pairs.

```
8978 --> 39
225 --> 60
89 --> 60
1715 --> 604
604 --> 3966
48 --> 1121
65 --> 1327
1715 --> 39
65 --> 19
2238 --> 39
32 --> 286
32 --> 156
156 --> 39
2238 --> 38
249 --> 36
855 --> 39
32 --> 101
352 --> 1859
9 --> 310
41 --> 1327
32 --> 271
48 --> 1715
1198 --> 39
749 --> 39
65 --> 9
48 --> 301
41 --> 310
41 --> 301
32 --> 310
41 --> 101
48 --> 1344
1344 --> 39
48 --> 1327
1327 --> 39
48 --> 271
65 --> 201
48 --> 1146
41 --> 1146
```

```
1146 --> 39
32 --> 110
208 --> 41
1121 --> 39
352 --> 39
249 --> 237
976 --> 41
41 --> 117
32 --> 117
522 --> 39
48 --> 9
48 --> 846
48 --> 740
19 --> 38
533 --> 39
48 --> 201
38 --> 271
48 --> 604
48 --> 110
41 --> 475
604 --> 39
48 --> 522
48 --> 310
48 --> 270
48 --> 225
41 --> 225
41 --> 170
41 --> 237
48 --> 475
475 --> 39
48 --> 438
438 --> 39
89 --> 41
65 --> 41
48 --> 60
60 --> 39
48 --> 352
41 --> 147
48 --> 161
310 --> 39
38 --> 286
48 --> 237
271 --> 39
270 --> 39
225 --> 39
48 --> 255
255 --> 39
48 --> 65
237 --> 39
36 --> 237
65 --> 39
179 --> 39
32 --> 179
32 --> 38
32 --> 48
32 --> 39
48 --> 170
```

```
170 --> 39
170 --> 38
48 --> 147
41 --> 110
110 --> 39
110 --> 38
89 --> 38
48 --> 101
48 --> 94
48 --> 89
101 --> 39
94 --> 39
89 --> 39
48 --> 79
41 --> 79
48 --> 41
48 --> 36
79 --> 39
32 --> 41
48 --> 39
48 --> 38
41 --> 39
41 --> 38
38 --> 39
37 --> 38
41 --> 36
36 --> 39
36 --> 38
None
```

with 20 percent of the dataset I got 117 pairs.

```
8978 --> 39
225 --> 60
89 --> 60
1715 --> 604
48 --> 1121
65 --> 1327
2238 --> 39
32 --> 156
156 --> 39
2238 --> 38
249 --> 36
855 --> 39
352 --> 1859
9 --> 310
89 --> 1327
41 --> 1327
32 --> 271
1198 --> 39
749 --> 39
65 --> 9
48 --> 301
41 --> 310
41 --> 301
32 --> 310
```

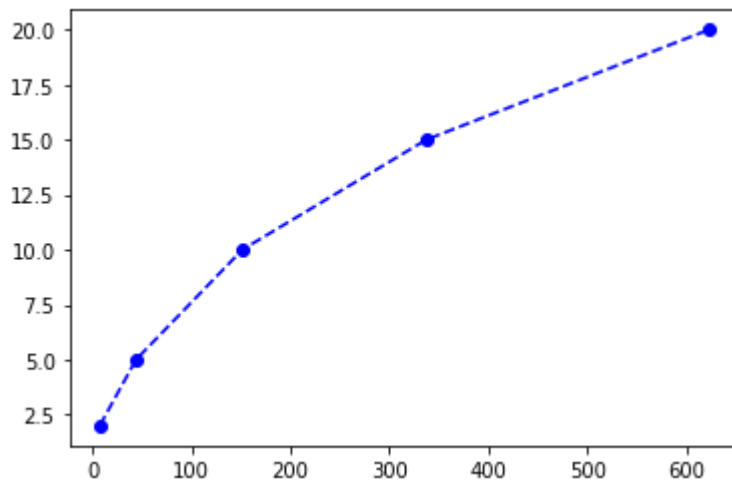
```
41 --> 101
48 --> 1344
1344 --> 39
48 --> 1327
1327 --> 39
48 --> 271
48 --> 117
65 --> 201
48 --> 1146
41 --> 1146
1146 --> 39
32 --> 110
208 --> 41
1121 --> 39
352 --> 39
249 --> 237
976 --> 41
41 --> 117
32 --> 117
41 --> 161
522 --> 39
48 --> 846
824 --> 39
769 --> 647
48 --> 740
533 --> 39
48 --> 201
38 --> 271
48 --> 604
48 --> 110
41 --> 475
604 --> 39
48 --> 310
48 --> 225
41 --> 225
41 --> 170
32 --> 170
41 --> 237
48 --> 475
475 --> 39
48 --> 438
438 --> 39
89 --> 41
65 --> 41
48 --> 60
60 --> 39
48 --> 352
41 --> 147
48 --> 161
310 --> 39
38 --> 286
48 --> 237
271 --> 39
270 --> 39
225 --> 39
48 --> 255
255 --> 39
```

```
48 --> 65
237 --> 39
36 --> 237
65 --> 39
179 --> 39
32 --> 179
32 --> 38
32 --> 48
32 --> 39
48 --> 170
170 --> 39
170 --> 38
48 --> 147
147 --> 39
41 --> 110
110 --> 39
110 --> 38
89 --> 38
48 --> 101
48 --> 89
101 --> 39
89 --> 39
48 --> 79
48 --> 41
48 --> 36
79 --> 39
32 --> 41
48 --> 39
48 --> 38
41 --> 39
41 --> 38
38 --> 39
37 --> 38
41 --> 36
36 --> 39
36 --> 38
None
```

## Graph

```
In [11]: 1 plt.plot(times, [x for x in sections], 'bo--')
```

```
Out[11]: [<matplotlib.lines.Line2D at 0x7f9130dd1640>]
```



## Netflix

```
In [ ]: 1 netflix = pd.read_csv('netflix.data', delimiter=" ", engine='python', on_bad_
2 netflix_L = netflix.values.tolist())
```

```
In [ ]: 1 netflix_NN = []
2 for x in netflix_L:
3     netflix_NN.append([i for i in x if str(i) != 'nan'])
```

```
In [ ]: 1 sections = [2,5,10,15, 20]
2 associationN = []
3 for i, x in enumerate(sections):
4     new_list = netflix_NN[0:int(len(netflix_NN)*(sections[i]/100))]
5     associationN.append(new_list)
```

```
In [ ]: 1 times = []
2 pairslen = []
3 actualpairs = []
4 for x in range(5):
5     print(x)
6     results = apyori(associationN[x], 0.01, 2000)
7     times.append(results[0])
8     pairslen.append(len(results[1]))
9     actualpairs.append(results[1])
```

```
In [ ]: 1 for x in range(len(sections)):
2     print()
3     print()
4     print('with %g percent of the dataset I got %d pairs.' % (sections[x], p
5     print()
6     print(output(actualpairs[x]))
```



```
In [ ]: 1 pl.plot(times, [x for x in sections], 'bo--')
```

## Conclusions

This is was the Multistage implementation of the Apyori algorithm. I feel that this was faster than the library algorithm as well as the normal apyori algorithm, reason being that I was able to run a max of 20% of the dataset and was not taking as long as the other algorithms. As for the pairs I also believed I maintained a consistent number of pairs throughout all my subsections. Unfortunately I am still not able to run my netflix data set. I've added a picture that you can see the contains the "kernel has died" error.