

Student name: Thien Toan Tran
Student ID: A1808080

Assignment 3 Exercise 1

In [22]:

```
1 from itertools import combinations
2 import random
3 import os
4 from datetime import datetime
5 import time
6 from multiprocessing import Pool
```

In [23]:

```
1 dataPath = './data/'
2 fileDataPath = []
3 fileDataNamesOnly = []
4 for a in os.listdir(dataPath):
5     if a.endswith('.dat'):
6         fileDataPath.append(dataPath + a)
7         fileDataNamesOnly.append(a.split('.')[0])
8
9 fileDataPath = sorted(fileDataPath)
10 fileDataNamesOnly = sorted(fileDataNamesOnly)
11 fileDataPath
```

Out[23]:

```
['./data/T10I4D100K.dat',
 './data/T40I10D100K.dat',
 './data/chess.dat',
 './data/connect.dat',
 './data/mushroom.dat',
 './data/pumsb.dat',
 './data/pumsb_star.dat']
```

In [24]:

```
1 dataLineCounts = dict()
2 for path in fileDataPath:
3     dataLineCounts[path] = sum(1 for line in open(path))
4 dataLineCounts
```

Out[24]:

```
{ './data/T10I4D100K.dat': 100000,
  './data/T40I10D100K.dat': 100000,
  './data/chess.dat': 3196,
  './data/connect.dat': 67557,
  './data/mushroom.dat': 8124,
  './data/pumsb.dat': 49046,
  './data/pumsb_star.dat': 49046}
```

Exercise 1: Frequent Itemsets

1.1. Implement the simple, randomized algorithm given in 6.4.1

Suppose need to randomly select n out of m baskets in the entire file, where $\frac{n}{m} = p$

My random algorithm will be implemented as follow:

1. select random subset of entire dataset
 - For each of basket in data file, a random float in the range $[0.0, 1.0)$ will be generated. It will represent the selection probability for this basket.
 - If the probability $\geq 1 - p$, then this basket will be selected for further process
2. Full itemset is initialized as a Python dictionary, where key is an itemset, value is its count on the subset
3. For each basket data, generate a set of itemset and update the full itemset
 - since the singleton and pairs are likely to dominate, only those 2 types are generated
 - for each generated itemset,
 - if it exists in the full itemset dictionary, then its value is increased by 1
 - otherwise, create new record using itself as key and value is 1
4. count itemset over the subset
5. sort itemsets and get most frequent itemset in the subset
 - lower the confident to ps , since the original data has expected confident s and we only process a portion p of the data
 - lower the confident more, like $0.9ps$, if there is enough memory. By doing that, the result is more likely to achieve the confident s in the whole dataset
 - equivalently change to least support count $0.9psm$

In [25]:



```

1 def randomSelect(m,p):
2     """
3     Randomly select ids based on given probability
4
5     :param m: dataset's record
6     :param p: probability
7     :return: list of selected ids
8     """
9     result = []
10    for i in range(m):
11        if random.random().__ge__(1-p):
12            result.append(i)
13    return result

```

Test random select result

In [26]:



```

1 list_m = [100,500,1000,10000]
2 list_p = [0.5,0.8,0.3]
3 for m in list_m:
4     for p in list_p:
5         a = randomSelect(m, p)
6         expect = m*p
7         selected = len(a)
8         pc = (expect- selected)/expect
9         print(f'm={m}, p={p}: expected: {expect}, result: {selected}, delta = {pc}')
10    print()
11

```

m=100, p=0.5: expected: 50.0, result: 51, delta = -2.0%

m=100, p=0.8: expected: 80.0, result: 71, delta = 11.25%

m=100, p=0.3: expected: 30.0, result: 31, delta = -3.3333333333333335%

m=500, p=0.5: expected: 250.0, result: 268, delta = -7.199999999999999%

m=500, p=0.8: expected: 400.0, result: 395, delta = 1.25%

m=500, p=0.3: expected: 150.0, result: 154, delta = -2.666666666666667%

m=1000, p=0.5: expected: 500.0, result: 486, delta = 2.8000000000000003%

m=1000, p=0.8: expected: 800.0, result: 790, delta = 1.25%

m=1000, p=0.3: expected: 300.0, result: 302, delta = -0.6666666666666667%

m=10000, p=0.5: expected: 5000.0, result: 4984, delta = 0.32%

m=10000, p=0.8: expected: 8000.0, result: 7947, delta = 0.6625%

m=10000, p=0.3: expected: 3000.0, result: 3017, delta = -0.5666666666666667%

In [27]:



```

1 # def testBasketIsSorted(basket: str):
2 #     item = basket.split('\s')
3 #     i = 0
4 #     while i < len(item)-2:
5 #         if item[i]>= item[i+1]: return False
6 #     return True

```

In [28]:



```
1 def generateItemset(basketItem:[int]):
2     """
3     Generate singleton and pairs from basket's item id
4
5     :param basketItem: array of basket's item id
6     :return: set of Itmeset
7     """
8     result = set()
9     for i in range(1, 4):
10         tmpset = combinations(basketItem,i)
11         for j in tmpset: result.add(j)
12
13     # i = 0
14     # while i < len(basketItem):
15     #     result.add((basketItem[i]))
16     #     j = i+1
17     #     while j < len(basketItem):
18     #         result.add((basketItem[i],basketItem[j]))
19     #         j+=1
20     #     i+=1
21     return result
```

In [29]:



```
1 def updateItemsetCount(fullItemset: dict, basketItemset):
2     """
3     Update basket's itemset into full itemset
4     if an itemset is already in the full itemset, increase its count by 1 in full
5     otherwise, add to the full itemset by creating new record of itself in full i
6
7     :param fullItemset: full itemset
8     :param basketItemset: basket itemset
9     :return:
10    """
11    for itemset in basketItemset:
12        if itemset in fullItemset:
13            # id = fullItemset.index(itemset)
14            fullItemset[itemset] +=1
15        else:
16            fullItemset[itemset] = 1
```

In [30]:



```
1 def convertLine(line:str):
2     """
3     Convert line into list of int
4
5     :param line: basket data
6     :return: list of item's id
7     """
8     id_str = line.split()
9     return list(map(int, id_str))
```

In [31]:



```
1 def processFileLine(line):
2     ids = convertLine(line)
3     return generateItemset(ids)
```

In [32]:



```
1 def filterResult(itemset:dict, leastSupportCount)-> dict:
2     """
3     Filter input itemset for only accepting ones with frequent greater or equal g
4     Input itemset must be sorted for correct result
5
6     :param itemset: collected itemsets and their frequent
7     :param leastSupportCount: given threshold
8     :return: itemsets that satisfy the condition
9     """
10    result = dict()
11    for key in itemset.keys():
12        if itemset[key] < leastSupportCount: break
13        result[key] = itemset[key]
14    return result
```

In [33]:



```

1  def processLimitedPassAlgo(filePath, m, p, s, resultLogPath, random, rangeId):
2      """
3      Modified version for run both randomize and SON
4
5      :param filePath:
6      :param m:
7      :param p: selecting ratio for randomize, 1/c for SON
8      :param s:
9      :param resultLogPath:
10     :param random: True for run randomized , False for run SON algo
11     :param rangeId: only need when run SON
12     :return: limited-pass frequent itemset
13     """
14
15     numberBasket = m
16
17     leastSupportCount = 0
18     intro = ''
19
20     processLineId = []
21     if random:
22         processLineId = randomSelect(numberBasket, p)
23         leastSupportCount = int(0.9*p*s*numberBasket)
24         intro = f'processing {filePath}, numberBasket={numberBasket}, numberSelec
25         intro += f', leastSupportCount={leastSupportCount}'
26     else:
27         processLineId = rangeId
28         leastSupportCount = int(s*numberBasket)
29         intro = f'processing {filePath}[{processLineId[0]}:{processLineId[-1]}]'
30     fullItemset = dict()
31
32     if random: print(intro)
33
34     timeStart = datetime.now()
35     currentLine = 0
36     with open(filePath, 'r') as file:
37         for line in file:
38             line = line.strip()
39             if (currentLine in processLineId) and (len(line)>0):
40                 updateItemsetCount(fullItemset, processFileLine(line))
41
42             currentLine += 1
43     # for itemset in fullItemset:
44     #     itemset.conf = itemset.count/numberBasket
45
46     frequentItemset = dict(sorted(fullItemset.items(),key=lambda item: item[1], r
47     if random:
48         frequentItemset = filterResult(frequentItemset,leastSupportCount)
49     timeStop = datetime.now()
50
51     end = f'finish in {str(timeStop - timeStart)}'
52     if random: print(end)
53
54     with open(resultLogPath,'w') as file:
55         file.write(f'{intro}\n{end}\n')
56         for key in frequentItemset.keys():
57             file.write(f'{key}->{frequentItemset[key]}\n')
58
59     return frequentItemset

```

In [34]:



```
1 def createDirectory(path):
2     """
3     Build folder based on given path
4     Support nested folder, ie. './a/b/c',
5     """
6     path = path.split(os.sep)
7     currentDir = '.'
8     for i in path:
9         currentDir += os.sep + i
10        try:
11            os.mkdir(currentDir)
12        except:
13            # this folder is existed
14            continue
15 p=.05
16 s=.5
17 dataDir = './data/'
18 resultDir = f'{dataDir}/p{p}s{s}/'
19 createDirectory(resultDir)
20
21 fileDataPath = []
22 for a in os.listdir(dataDir):
23     if a.endswith('.dat'): fileDataPath.append(dataDir + a)
24 fileDataPath = sorted(fileDataPath)
```

Run randomised algo as requested in 1.3

Dataset for running need to be found in folder "data"

Result will be saved in folder with format "data/p{s}s{s}", ex. "data/p0.05s0.5/"

In [14]:



```

1 for path in fileDataPath:
2     m = sum(1 for line in open(path))
3     outputPath = path.replace('.dat', '.resultv2')\
4                 .replace(dataDir, resultDir)
5     processLimitedPassAlgo(path, m, p, s, outputPath, True, [])
6

```

```

processing ./data/T10I4D100K.dat, numberBasket=100000, numberSelecting=
5120, leastSupportCount=2250
finish in 0:00:06.740719
processing ./data/T40I10D100K.dat, numberBasket=100000, numberSelecting
=4973, leastSupportCount=2250
finish in 0:01:43.280514
processing ./data/chess.dat, numberBasket=3196, numberSelecting=153, le
astSupportCount=71
finish in 0:00:00.583342
processing ./data/connect.dat, numberBasket=67557, numberSelecting=344
7, leastSupportCount=1520
finish in 0:00:31.512795
processing ./data/mushroom.dat, numberBasket=8124, numberSelecting=398,
leastSupportCount=182
finish in 0:00:00.295855
processing ./data/pumsb.dat, numberBasket=49046, numberSelecting=2468,
leastSupportCount=1103
finish in 0:02:57.801484
processing ./data/pumsb_star.dat, numberBasket=49046, numberSelecting=2
393, leastSupportCount=1103
finish in 0:00:52.036848

```

In []:



1

1.2. Implement SON algorithm

SON algorithm is implemented as follow phases: **Phase 1:**

1. Divide full dataset into c chunks
 - modify RandomAlgo to take a range of line Id instead of randomly generating
 - equivalently $p = \frac{1}{c}$
2. For each of data chunk:
 - handled by a parallel 'process'
 - each 'process' execute same process as Random Algorithm
 - store outputs into files

Phase 2:

1. Collect and summarize frequent itemset from output files
2. Sort itemsets and get most frequent itemset
 - expected confident is s
 - equivalently change to least support count sm

In []:



1

Implementation for phase 1

I chose *multiprocessing.Pool* for implementing parallel processing.

Pool will be initialized with a specific amount of subprocess, ie. 5 in my implementation. Those subprocess is automatically mapped with given function and given arguments set, then will be executed independently from the main process

In [35]:



```

1 def buildDataForPool(path, m, p, s, outputPath):
2     """
3     Build list of arguments list for pool's subprocess
4
5     :param filePath: path to dataset
6     :param m:
7     :param p:
8     :param s:
9     :param outputPath: output path for chunk's result
10    :return: list of arguments list for pool's subprocess
11    """
12
13    data = []
14    mOverC = int(m/c)
15    maxCounter = c if m%c==0 else c+1
16
17    for counter in range(maxCounter):
18        # each item must strictly follow the argument order from
19        # processLimitedPassAlgo(filePath, m, p, s, resultLogPath, random, rangeId)
20        data.append([path, m, p, s, outputPath+str(counter), False
21                    , [x for x in range(mOverC*counter, min(m, mOverC*(counter
22                    ))]
23    return data
24
25

```

In [36]:



```

1 def parallelProcess(processInput):
2     """
3     Execute processLimitedPassAlgo with given argument list
4     """
5
6     # processLimitedPassAlgo(filePath, m, p, s, resultLogPath, random, rangeId):
7     processLimitedPassAlgo(processInput[0], processInput[1], processInput[2], proces
8                             , processInput[4], processInput[5], processInput[6])
9
10

```

In [17]:



```

1 pool = Pool(5)
2 c = 100
3 resultDir = f'{dataDir}SONc{c}s{s}/'
4 createDirectory(resultDir)

```

Implementation for phase 2:

In [18]:



```

1 def collectItemsetFromChunks(chunkResultList):
2     """
3     Combine full itemset result from all created results
4
5     :param chunkResultList: list of chunk's result files, ie. chess.result.sonXX
6     :return: full itemset
7     """
8
9     result = dict()
10    for filePath in chunkResultList:
11        # print(f'executing collectItemsetFromChunks: {filePath}')
12        with open(filePath, 'r') as file:
13            # skip 2 info lines
14            next(file)
15            next(file)
16
17            for line in file:
18                key, value = line.split('->')
19                value = int(value)
20                if key in result: result[key] += int(value)
21                else: result[key] = int(value)
22    return result
23

```

In [19]:



```

1 def collectItemsetFromDataset(datasetName, resultDir):
2     """
3     SON phase 2's functionalities
4     """
5
6     timeStart = datetime.now()
7     m = dataLineCounts[f'{dataDir}{datasetName}.dat']
8     chunkResult = datasetName + '.result.son'
9
10    # scan resultDir for related chunk result
11    chunkResultList = [resultDir + x for x in os.listdir(resultDir) if chunkResult
12
13
14    fullItemset = collectItemsetFromChunks(chunkResultList)
15    fullItemset = dict(sorted(fullItemset.items(), key=lambda item: item[1], reverse=True))
16    timeStop = datetime.now()
17
18    print(f'finish collect from {chunkResult} in {str(timeStop - timeStart)}')
19
20    with open(f'{resultDir}{datasetName}.final', 'w') as file:
21        for key in fullItemset.keys():
22            if fullItemset[key] < s*m: break
23            file.write(f'{key}->{fullItemset[key]}\n')
24    # return fullItemset

```

In []:



1

Run SON algorithm as request in 1.3

In [20]:



```

1 %%time
2
3 # run SON's phase 1
4
5 for path in fileDataPath:
6     timeStart = datetime.now()
7     m = dataLineCounts[path]
8     outputPath = path.replace('.dat', '.result.son')\
9                     .replace(dataDir, resultDir)
10
11     parallelData = buildDataForPool(path, m, 1/c, s, outputPath)
12     pool.map(parallelProcess, parallelData)
13     timeStop = datetime.now()
14     print(f'finsih processing {path} in {str(timeStop-timeStart)}')
```

```

finsih processing ./data/T10I4D100K.dat in 0:01:12.780084
finsih processing ./data/T40I10D100K.dat in 0:15:20.015748
finsih processing ./data/chess.dat in 0:00:06.403641
finsih processing ./data/connect.dat in 0:04:33.407238
finsih processing ./data/mushroom.dat in 0:00:03.636486
finsih processing ./data/pumsb.dat in 0:23:17.710465
finsih processing ./data/pumsb_star.dat in 0:07:42.315305
CPU times: user 195 ms, sys: 24.8 ms, total: 219 ms
Wall time: 52min 16s
```

In [21]:



```

1 %%time
2
3 # run SON's phase 2
4
5 for filename in fileDataNamesOnly:
6     collectItemsetFromDataset(filename, resultDir)
7
```

```

finish collect from T10I4D100K.result.son in 0:00:26.140894
finish collect from T40I10D100K.result.son in 0:17:00.326199
finish collect from chess.result.son in 0:00:01.601326
finish collect from connect.result.son in 0:00:07.835119
finish collect from mushroom.result.son in 0:00:01.119917
finish collect from pumsb.result.son in 0:07:01.125634
finish collect from pumsb_star.result.son in 0:03:28.021837
CPU times: user 27min 28s, sys: 34.8 s, total: 28min 3s
Wall time: 28min 12s
```

In []:



1

1.3. Report the outcomes

a3e13_result.zip contains the outcomes of procssing following datasets:

- T10I4D100K.dat
- T40I10D100K.dat

- chess.dat
- connect.dat
- mushroom.dat
- pumsb.dat
- pumsb_star.dat

Folder includes:

- folder **SONc100s0.5** contains only final results of SON's algo running with 100 chunks and s=0.5
- folder **p0.05s0.5** contains result of randomised algo running with p=0.05 and s=0.5

1.4. Experiment with different sample sizes in the simple randomized algo-rithm such as 1, 2, 5, 10% and compare your results

In [41]:

```

1 p_list = [.01,.02,.1,.2]
2
3 for p in p_list:
4     s=.5
5     dataDir = './data/'
6     resultDir = f'{dataDir}/p{p}s{s}/'
7     createDirectory(resultDir)
8
9     for path in fileDataPath:
10        m = sum(1 for line in open(path))
11        outputPath = path.replace('.dat','.resultv2')\
12                      .replace(dataDir,resultDir)
13        processLimitedPassAlgo(path, m, p, s, outputPath, True, [])
14    print()
```

```

processing ./data/T10I4D100K.dat, numberBasket=100000, numberSelecting
=998, leastSupportCount=450
finish in 0:00:01.349564
processing ./data/T40I10D100K.dat, numberBasket=100000, numberSelectin
g=994, leastSupportCount=450
finish in 0:00:13.346378
processing ./data/chess.dat, numberBasket=3196, numberSelecting=31, le
astSupportCount=14
finish in 0:00:00.139196
processing ./data/connect.dat, numberBasket=67557, numberSelecting=71
2, leastSupportCount=304
finish in 0:00:05.447991
processing ./data/mushroom.dat, numberBasket=8124, numberSelecting=10
5, leastSupportCount=36
finish in 0:00:00.094768
processing ./data/pumsb.dat, numberBasket=49046, numberSelecting=499,
leastSupportCount=220
finish in 0:00:31.899542
processing ./data/pumsb_star.dat, numberBasket=49046, numberSelecting=
```

In []:

1

In []:



1

T10I4D100K.dat and T40I10D100K.dat

Randomised algo's results and SON algo's result both show empty itemset. We can conclude that itemsets in those dataset are fragmented, therefore they don't have enough repetition to achieve expected confidence s

chess.dat's top 5 frequent item

p=1% s=0.5	p=2% s=0.5	p=5% s=0.5	p=10% s=0.5	p=20% s=0.5	SON c=100 s=0.5
(48, 58)->31	(58,)->66	(52, 58)->153	(58,)->310	(58,)->625	(58,)->3130
(40, 48, 58)->31	(60,)->66	(58,)->153	(52, 58)->309	(29,)->624	(60,)->3127
(52, 60)->31	(58, 60)->66	(52,)->153	(52,)->309	(52, 58)->624	(58, 60)->3126
(40, 48, 60)->31	(29, 58, 60)->65	(40,)->152	(29,)->308	(29, 58)->624	(52,)->3120
(40,)->31	(52, 60)->65	(40, 52, 58)->152	(29, 58)->308	(52,)->624	(52, 58)->3119

Consider as SON has the most accurate (process 100% dataset instead of picking random parts)

- p=20%'s result contains 3 correct frequent itemset
- p=10%'s result contains 2 correct frequent itemset
- p=5%'s result contains 3 correct frequent itemset
- p=2%'s result contains 1 correct frequent itemset
- p=1%'s result contains 0 correct frequent itemset

The incorrect itemset appear in other results are called *false negative*

The higher value of p , the lesser false negative we get.

The frequent of itemset seems to be more accurate with higher p .

- SON result indicates singleton (58) repeats 3130 times
- p=20%'s indicates singleton (58) repeats 625 times -> $625 \times 5 = 3125$ times over the whole dataset: extremely closed to SON's result
- p=10%'s indicates singleton (58) repeats 310 times -> $310 \times 10 = 3100$ times over the whole dataset: closed
- p=5%'s indicates singleton (58) repeats 153 times -> $153 \times 20 = 3060$ times over the whole dataset: the difference is ~100
- p=2%'s indicates singleton (58) repeats 66 times -> $66 \times 50 = 3300$ times over the whole dataset: the difference is ~170
- p=1%'s indicates singleton (58) repeats 31 times -> $31 \times 100 = 3100$ times over the whole dataset: really closed, but (58) is only the 19th among most frequent itemset, while it actually the 1st one in final

Similar analyze could be done based other dataset's results