

# CS145 Howework 5

**\*\*Important Note:\*\*** HW4 is due on **11:59 PM PT, Dec 4 (Friday, Week 9)**. Please submit through GradeScope.

## Print Out Your Name and UID

**\*\*Name:** Rui Deng, **UID:** 205123245**\*\***

## Before You Start

You need to first create HW5 conda environment by the given `cs145hw5.yml` file, which provides the name and necessary packages for this tasks. If you have `conda` properly installed, you may create, activate or deactivate by the following commands:

```
conda env create -f cs145hw5.yml
conda activate hw4
conda deactivate
```

OR

```
conda env create --name NAMEOFOURCHOICE -f cs145hw5.yml
conda activate NAMEOFOURCHOICE
conda deactivate
```

To view the list of your environments, use the following command:

```
conda env list
```

More useful information about managing environments can be found [here](https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html) (<https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html>).

You may also quickly review the usage of basic Python and Numpy package, if needed in coding for matrix operations.

In this notebook, you must not delete any code cells in this notebook. If you change any code outside the blocks (such as some important hyperparameters) that you are allowed to edit (between START/END YOUR CODE HERE), you need to highlight these changes. You may add some additional cells to help explain your results and observations.

In [1]:

```
import numpy as np
import pandas as pd
import sys
import random
import math
import matplotlib.pyplot as plt
from graphviz import Digraph
from IPython.display import Image
from scipy.stats import multivariate_normal
%load_ext autoreload
%autoreload 2
```

If you can successfully run the code above, there will be no problem for environment setting.

## 1. Frequent Pattern Mining for Set Data (25 pts)

**Table 1**

TID	Items
1	b,c,j
2	a,b,d
3	a,c
4	b,d
5	a,b,c,e
6	b,c,k
7	a,c
8	a,b,e,i
9	b,d
10	a,b,c,d

Given a transaction database shown in Table 1, answer the following questions. Let the parameter `min_support` be 2.

### Questions

#### 1.1 Apriori Algorithm (16 pts) .

**Note:** This is a "question-answer" style problem. You do not need to code anything and you are required to calculate by hand (with a scientific calculator). Find all the frequent patterns using Apriori Algorithm.

- $C_1$
- $L_1$
- $C_2$
- $L_2$
- $C_3$
- $L_3$
- $C_4$
- $L_4$

Please type your answer here!

a.  $C_1 = \{a, b, c, d, e, i, j, k\}$

b.  $L_1 = \{a, b, c, d, e\}$ . We scan the database and find out i, j, k only have support = 1 < min\_support, while others all have support  $\geq 2$ .

c.  $C_2 = \{ab, ac, ad, ae, bc, bd, be, cd, ce, de\}$

d.  $L_2 = \{ab, ac, ad, ae, bc, bd, be\}$ . We scan the database and find out cd and ce only have support = 1 < min\_support, and de has support = 0 < min\_support; while others all have support  $\geq 2$ .

e. After Self joining:  $C'_3 = \{abc, abd, abe, acd, ace, ade, bcd, bce, bde\}$

After Pruning:  $C_3 = \{abc, abd, abe, acd, ace\}$ . Here,  $\{ade, bcd, bce, bde\}$  can be safely pruned because de, cd, ce are not in  $L_2$ .

f.  $L_3 = \{abc, abd, abe\}$ . We scan the database and find out acd and ace only have support = 1 < min\_support; while others all have support  $\geq 2$ .

g. After Self joining:  $C'_4 = \{abcd, abce, abde\}$

After Pruning:  $C_4 = \{\}$ . Here,  $\{abcd, abce, abde\}$  can be safely pruned because acd, ace, ade are not in  $L_3$ .

h.  $L_4 = \{\}$

So the frequent patterns are  $\{a, b, c, d, e, ab, ac, ad, ae, bc, bd, be, abc, abd, abe\}$ .

## 1.2 FP-tree (9 pts)

(a) Construct the FP-tree of the table.

(b) For the item d, show its conditional pattern base (projected database) and conditional FP-tree

You may use Package `graphviz` to generate graph

(<https://graphviz.readthedocs.io/en/stable/manual.html>

(<https://graphviz.readthedocs.io/en/stable/manual.html>)) (Bonus point: 5pts) or draw by hand.

(c) Find frequent patterns based on d's conditional FP-tree

Please type your answer here!

(a)

In [28]:

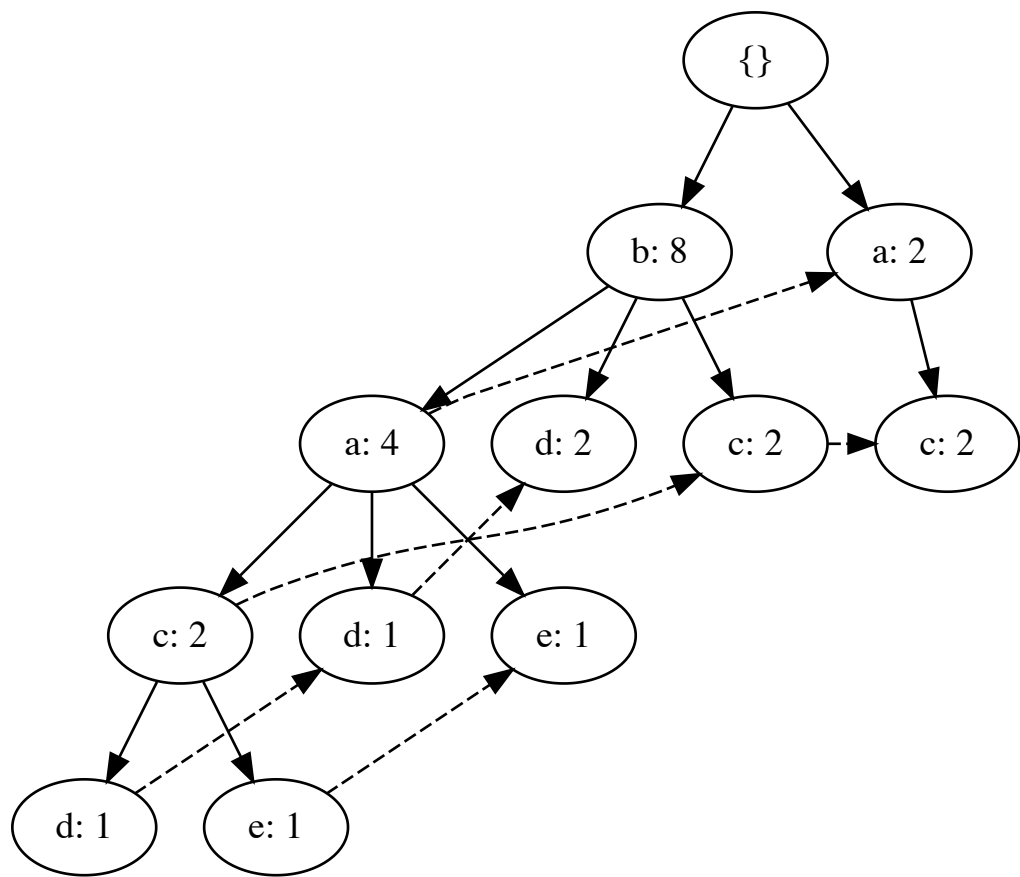
```
d = Digraph()
d.node('1', '{}')
d.node('2', 'b: 8')
d.node('3', 'a: 2')
d.node('4', 'a: 4')
d.node('5', 'd: 2')
d.node('6', 'c: 2')
d.node('7', 'c: 2')
d.node('8', 'c: 2')
d.node('9', 'd: 1')
d.node('10', 'e: 1')
d.node('11', 'd: 1')
d.node('12', 'e: 1')

d.edge('1', '2')
d.edge('1', '3')
d.edge('3', '7')
d.edge('2', '4')
d.edge('2', '5')
d.edge('2', '6')
d.edge('4', '8')
d.edge('4', '9')
d.edge('4', '10')
d.edge('8', '11')
d.edge('8', '12')

d.edge('4', '3', style="dashed", constraint="false")
d.edge('8', '6', style="dashed", constraint="false")
d.edge('6', '7', style="dashed", constraint="false")
d.edge('11', '9', style="dashed", constraint="false")
d.edge('9', '5', style="dashed", constraint="false")
d.edge('12', '10', style="dashed", constraint="false")

d
```

Out[28]:



(b)

Conditional pattern base for d is: {abc: 1, ab: 1, b: 2}. Then we construct its conditional FP-tree:

In [31]:

```

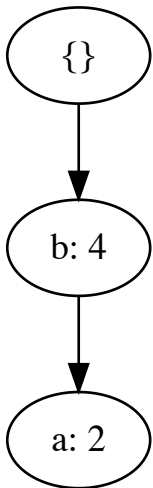
d_b = Digraph()
d_b.node('1', '{}')
d_b.node('2', 'b: 4')
d_b.node('3', 'a: 2')

d_b.edge('1', '2')
d_b.edge('2', '3')

d_b

```

Out[31]:



(c)

So the frequent patterns for d is: {abd, ad, bd, d}

## 2. Apriori for Yelp (50 pts)

In `apriori.py`, fill the missing lines. The parameters are set as `min_support=50` and `min_conf = 0.25`, and `ignore_one_iter_set=True`. Use the Yelp data `yelp.csv` and `id_nams.csv`, and run the following cell and report the frequent patterns and rules associated with it.

In [35]:

```
#No need to modify
from hw5code.apriori import *
input_file = read_data('./data/yelp.csv')
min_support = 50
min_conf = 0.25
items, rules = run_apriori(input_file, min_support, min_conf)
name_map = read_name_map('./data/id_name.csv')
print_items_rules(items, rules, ignore_one_item_set=True, name_map=name_map)
```

```
item:
"Wicked Spoon","Holsteins Shakes & Buns" 51
item:
"Wicked Spoon","Earl of Sandwich" 52
item:
"Wicked Spoon","Secret Pizza" 52
item:
"Wicked Spoon","The Cosmopolitan of Las Vegas" 54
item:
"Mon Ami Gabi","Wicked Spoon" 57
item:
"Bacchanal Buffet","Wicked Spoon" 63

----- RULES:
Rule:
"Secret Pizza" "Wicked Spoon" 0.2561576354679803
Rule:
"The Cosmopolitan of Las Vegas" "Wicked Spoon" 0.27692307692307694
Rule:
"Holsteins Shakes & Buns" "Wicked Spoon" 0.3148148148148148
```

What do these results mean? Do a quick Google search and briefly interpret the patterns and rules mined from Yelp in 50 words or less.

[Please type your answer here!](#)

The rules mean that these businesses are highly associated with each other. The Cosmopolitan of Las Vegas is a luxury resort, and the Secret Pizza, Wicked Spoon, and Holsteins Shakes & Buns are restaurants nearby. The high association means that people usually visit these place together (for example, they live at the resort and take meals in these restaurants).

### 3. Correlation Analysis (10 pts)

Note: This is a "question-answer" style problem. You do not need to code anything and you are required to calculate by hand (with a scientific calculator).

**Table 2**

---	Beer	No Beer	Total
Nuts	150	700	850
No Nuts	350	8800	9150
Total	500	9500	10000

Table 2 shows how many transactions containing beer and/or nuts among 10000 transactions.

Answer the following questions:

3.1 Calculate `confidence`, `lift` and `all_confidence` between buying beer and buying nuts.

3.2 What are your conclusions of the relationship between buying beer and buying nuts? Justify your conclusion with the previous measurements you calculated in 3.1.

Please type your answer here!

1.

$$\begin{aligned} \text{confidence}(\text{Beer} \rightarrow \text{Nuts}) &= \frac{\text{support}(\text{BeerAndNuts})}{\text{support}(\text{Beer})} = \frac{150}{500} = 0.3 \\ \text{confidence}(\text{Nuts} \rightarrow \text{Beer}) &= \frac{\text{support}(\text{BeerAndNuts})}{\text{support}(\text{Nuts})} = \frac{150}{850} = 0.17647 \\ \text{lift} &= \frac{P(\text{BeerAndNuts})}{P(\text{Beer})P(\text{Nuts})} = \frac{150}{850 * 500/10000} = 3.529 \\ \text{all\_confidence} &= \min(\text{confidence}(\text{Beer} \rightarrow \text{Nuts}), \text{confidence}(\text{Nuts} \rightarrow \text{Beer})) = 0.17647 \end{aligned}$$

1. Here we have  $\text{lift} > 1$ , which means that buying nuts and buying beer are positively correlated. Still, lift is not null-invariant, and its large positive value may be caused by the high number of transactions that do not contain beer or nuts. Then, we look at `all_confidence`, which is a null-invariant measure. In this case, we know that both transaction will appear together with a confidence of 0.176.

## 4. Sequential Pattern Mining (GSP Algorithm) (15 pts)

Note: This is a "question-answer" style problem. You do not need to code anything and you are required to calculate by hand (with a scientific calculator).

4.1 For a sequence  $s = \langle ab(cd)(ef) \rangle$ , how many events or elements does it contain? What is the length of  $s$ ? How many non-empty subsequences does  $s$  contain?

4.2 Suppose we have  $L_3 = \{ \langle (ac)e \rangle, \langle b(cd) \rangle, \langle bce \rangle, \langle a(cd) \rangle, \langle (ab)d \rangle, \langle (ab)c \rangle \}$ , as the frequent 3-sequences, write down all the candidate 4-sequences  $C_4$  with the details of the join and pruning steps.

Please type your answer here!

1. This sequence contains 4 events: {a, b, (cd), (ef)}, and its length is 6. The number of non-empty subsequences is  $2^6 - 1 = 63$



### 1. Joining:

Here, we can join  $\langle (ab)c \rangle$  and  $\langle b(cd) \rangle$  since they share common subsequence  $bc$ , and we get  $\langle (ab)(cd) \rangle$ .

We can also join  $\langle (ab)c \rangle$  with  $\langle bce \rangle$  since they also share the common subsequence  $bc$ , and we get  $\langle (ab)ce \rangle$ .

### Pruning:

For  $\langle (ab)(cd) \rangle$ , its subsequences of length 3 are  $\langle (ab)c \rangle$ ,  $\langle (ab)d \rangle$ ,  $\langle a(cd) \rangle$ ,  $\langle b(cd) \rangle$ , and they are all contained in  $L_3$ .

For  $\langle (ab)ce \rangle$ , its subsequences of length 3 are  $\langle (ab)c \rangle$ ,  $\langle (ab)e \rangle$ ,  $\langle ace \rangle$ ,  $\langle bce \rangle$ , and  $\langle (ab)e \rangle$ ,  $\langle ace \rangle$  are not contained in  $L_3$ , so we pruned this sequence.

Then,  $C_4 = \langle (ab)(cd) \rangle$

## 5 Bonus Question (10 pts)

1. In FP-tree, what will happen if we use ascending instead descending in header table?
2. Describe CloSpan (Mining closed sequential patterns: CloSpan (Yan, Han & Afshar @SDM'03)). Compare with algorithms we discussed in class.

Please type your answer here!

1. If we use ascending order in header table, then the less frequent items will appear before the more frequent items. This will cause the tree to grow wider and have more branches, since it is very possible that the frequent items are not included and thus not appear in the same branch with the top entries of the header table. Thus, the algorithm will take longer to run and use more spaces to store the FP-trees, which is an extra cost.
2. First, we know that a closed pattern is a pattern such that it is not included in another pattern having the same support. Then, the CloSpan aims to discover all frequent (i.e. with support  $\geq \text{min\_support}$ ) closed patterns in a database. The CloSpan is a further development of PrefixSpan that we learnt in class, since it prunes redundant patterns with Backward Subpattern and Backward Superpattern pruning and also attains the same information with the Prefix Span algorithm. It can efficiently reduce the search space and speeds up the pattern finding process.

## End of Homework 5 :)

After you've finished the homework, please print out the entire `ipynb` notebook and four `py` files into one PDF file. Make sure you include the output of code cells and answers for questions. Prepare submit it to GradeScope. Also this time remember assign the pages to the questions on GradeScope

```

from itertools import chain, combinations, islice
from collections import defaultdict
from time import time
import pandas as pd
import operator

def run_apriori(infile, min_support, min_conf):
    """
    Run the Apriori algorithm. infile is a record iterator.
    Return:
        rtn_items: list of (set, support)
        rtn_rules: list of ((preset, postset), confidence)
    """
    one_cand_set, all_transactions = gen_one_item_cand_set(infile)

    set_count_map = defaultdict(int) # maintains the count for each set

    one_freq_set, set_count_map = get_items_with_min_support(
        one_cand_set, all_transactions, min_support, set_count_map)

    freq_map, set_count_map = run_apriori_loops(
        one_freq_set, set_count_map, all_transactions, min_support)

    rtn_items = get_frequent_items(set_count_map, freq_map)
    rtn_rules = get_frequent_rules(set_count_map, freq_map, min_conf)

    return rtn_items, rtn_rules

def gen_one_item_cand_set(input_fileator):
    """
    Generate the 1-item candidate sets and a list of all the transactions.
    """
    all_transactions = list()
    one_cand_set = set()
    for record in input_fileator:
        transaction = frozenset(record)
        all_transactions.append(transaction)
        #=====#
        # START YOUR CODE HERE #
        #=====#
        for item in transaction:
            itemS = set()
            itemS.add(item)
            if frozenset(itemS) not in one_cand_set:
                one_cand_set.add(frozenset(itemS))
        #=====#
        # END YOUR CODE HERE #
        #=====#
    return one_cand_set, all_transactions

```

```

def get_items_with_min_support(item_set, all_transactions, min_support,
                               set_count_map):
    """
    item_set is a set of candidate sets.
    Return a subset of the item_set
    whose elements satisfy the minimum support.
    Update set_count_map.
    """
    rtn = set()
    local_set = defaultdict(int)

    for item in item_set:
        for transaction in all_transactions:
            if item.issubset(transaction):
                set_count_map[item] += 1
                local_set[item] += 1

    #=====#
    # STRART YOUR CODE HERE #
    #=====#
    for item, count in local_set.items():
        if count >= min_support:
            rtn.add(item)

    #=====#
    #  END YOUR CODE HERE  #
    #=====#

    return rtn, set_count_map


def run_apriori_loops(one_cand_set, set_count_map, all_transactions,
                       min_support):
    """
    Return:
        freq_map: a dict
                  {<length_of_set_l>: <set_of_frequent_itemsets_of_length_l>}
        set_count_map: updated set_count_map
    """
    freq_map = dict()
    current_l_set = one_cand_set
    i = 1
    #=====#
    # STRART YOUR CODE HERE #
    #=====#
    while (current_l_set != set([])):
        freq_map[i] = current_l_set
        current_l_set = join_set(current_l_set, i + 1)

```

```

        current_c_set, set_count_map =
            get_items_with_min_support(current_l_set, all_transactions,
                                      min_support, set_count_map)
        current_l_set = current_c_set

        i += 1
        #=====#
        #   END YOUR CODE HERE   #
        #=====#

    return freq_map, set_count_map

def get_frequent_items(set_count_map, freq_map):
    """ Return frequent items as a list. """
    rtn_items = []
    for key, value in freq_map.items():
        rtn_items.extend(
            [(tuple(item), get_support(set_count_map, item))
             for item in value])
    return rtn_items

def get_frequent_rules(set_count_map, freq_map, min_conf):
    """ Return frequent rules as a list. """
    rtn_rules = []
    for key, value in islice(freq_map.items(), 1, None):
        for item in value:
            _subsets = map(frozenset, [x for x in subsets(item)])
            for element in _subsets:
                remain = item.difference(element)
                if len(remain) > 0:
                    #=====#
                    #   STRART YOUR CODE HERE   #
                    #=====#
                    confidence = float(set_count_map[element.union(remain)]) /
                                float(set_count_map[element])
                    #=====#
                    #   END YOUR CODE HERE   #
                    #=====#
                    if confidence >= min_conf:
                        rtn_rules.append(
                            ((tuple(element), tuple(remain)), confidence))
    return rtn_rules

def get_support(set_count_map, item):
    """ Return the support of an item. """
    #=====#
    #   STRART YOUR CODE HERE   #
    #=====#

```

```

sup_item = set_count_map[item]
#=====#
#   END YOUR CODE HERE   #
#=====#
return sup_item

def join_set(s, l):
    """
    Join a set with itself .
    Return a set whose elements are unions of sets in s with length==l.
    """
    #=====#
    # STRART YOUR CODE HERE #
    #=====#
    join_set = set()
    for s1 in s:
        for ss1 in s:
            if s1 is not ss1:
                join = s1.union(ss1)

                if len(join) == l:
                    join_set.add(join)
    #   END YOUR CODE HERE   #
    #=====#
    return join_set

def subsets(x):
    """ Return non --empty subsets of x. """
    return chain(*[combinations(x, i + 1) for i, a in enumerate(x)])

def print_items_rules(items, rules, ignore_one_item_set=False, name_map=None):
    for item, support in sorted(items, key=operator.itemgetter(1)):
        if len(item) == 1 and ignore_one_item_set:
            continue
        print ('item: ')
        print (convert_item_to_name(item, name_map), support)
    print ('\n----- RULES:')
    for rule, confidence in sorted(
        rules, key=operator.itemgetter(1)):
        pre, post = rule
        print ('Rule: ')
        print( convert_item_to_name(pre, name_map), convert_item_to_name(post,
            name_map), confidence)

def convert_item_to_name(item, name_map):
    """ Return the string representation of the item. """
    if name_map:

```

```

        return ','.join([name_map[x] for x in item])
    else:
        return str(item)

def read_data(fname):
    """ Read from the file and yield a generator. """
    file_iter = open(fname, 'rU')
    for line in file_iter:
        line = line.strip().rstrip(',')
        record = frozenset(line.split(','))
        yield record

def read_name_map(fname):
    """ Read from the file and return a dict mapping ids to names. """
    df = pd.read_csv(fname, sep=',\t ', header=None, names=['id', 'name'],
                     engine='python')
    return df.set_index('id')['name'].to_dict()

```