

# UCLA CS 145 Homework #5

DUE DATE: Friday, 11/30/2018 11:59 PM

## Note

- You are expected to submit both a report and code. The submission format is specified on CCLE under HW5 description.
- Copying and sharing of homework are NOT allowed. But you can discuss general challenges and ideas with others. *Suspicious cases will be reported to The Office of the Dean of Students.*
- “# ===== YOUR CODE HERE =====”  
is used where input from you is needed in the code file.

## 1 Frequent Pattern Mining for Set Data

Given a transaction database shown in Table 1, answer the following questions. Note that the parameter `min_support` is set as 2.

- Find all the frequent patterns using Apriori Algorithm. Details of the procedure are expected.
- Construct and draw the FP-tree of the transaction database.
- For the item  $d$ , show its conditional pattern base (projected database) and conditional FP-tree.
- Find frequent patterns based on  $d$ 's conditional FP-tree.

Table 1: The transaction database for the question 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	$a, b, c, d$

## 2 Apriori for Yelp

In `apriori.py`, fill in the missing lines, with the following parameters (already set in the code): `min_support=50`, `min_conf=0.25`, and `ignore_one_item_set=True`. Output the frequent patterns and rules associated with the Yelp data (the same one as the project) which we have stored in `yelp.csv` and `id_name.csv`. Do NOT modify the `print_items_rules()` function and directly copy the entire output of the following command in your report in plain text format (do NOT take a screenshot):

```
python2.7 apriori.py
```

What patterns and rules do you see? Where are these businesses located? 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.

## 3 Correlation Analysis

Table 2 shows how many transactions containing beer and/or nuts among 10000 transactions. Answer the following questions based on Table 2.

- Calculate **confidence**, **lift**, and **all\_confidence** between buying beer and buying nuts.
- What are your conclusions of the relationship between buying beer and buying nuts, based on the above measures?

Table 2: Contingency table for question 2.

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

## 4 Sequential Pattern Mining (GSP Algorithm)

- 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?
- 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.