INFS 5102 - Unsupervised Methods in Analytics

Practical #10: Association Analysis in Python

Objective:

• Learn how to do association analysis using Python.

Submission:

- What to submit: The saved document (.html) containing the required steps done in this practical.
- <u>Deadline of the submission</u>: 11:59PM (Adelaide Time), Tuesday of Week 13.
- <u>Submission link</u>: "Submission Link of Prac #10" in Week 12 section on Learnonline course site.
- *Marks*: Prac#10 (part of the ongoing assessment of the course) is worth 2% of the total marks of the course.

Association Analysis in Python:

Association analysis is a rule-based unsupervised learning method for discovering interesting relations between variables in databases. The association rules can be determined by 3 parameters (support, confidence, and lift) that are used to identify the algorithm's strength. The **Apriori** algorithm is the most popular association rules approach, which can help us to find strong relationships between different items/products in the industry, such as retail industry data in this practical ("**Prac#10-GroceryStore.csv**").

The first step, as always, is to import the required Python libraries. Please execute the following script to do so:

```
import pandas as pd
import numpy as np
from apyori import apriori
from mlxtend.frequent_patterns import apriori, association_rules
```

Note: You must download and install **apyori** library in the default path for your Python libraries before proceeding.

Open the terminal and type in the following command, as shown in the screenshot below,

pip install apyori

```
Last login: Thu Jul 21 12:48:53 on ttys001

The default interactive shell is now zsh.
To update your account to use zsh, please run `chsh -s /bin/zsh`.
For more details, please visit https://support.apple.com/kb/HT208050.
[(base) STM285560:~ caozeh$ pip install apyori

Collecting apyori

Downloading apyori-1.1.2.tar.gz (8.6 kB)

Building wheels for collected packages: apyori

Building wheel for apyori (setup.py) ... done

Created wheel for apyori: filename=apyori-1.1.2-py3-none-any.whl size=5974 sha
256-0348cffda38c2f0f198e6627957e6cf7177f37b51d5aa1356fc114aaca7b9060

Stored in directory: /Users/caozeh/Library/Caches/pip/wheels/32/2a/54/10c59551
5f38bf3726642b10c60bf788029e8f3a1323e3913a
Successfully built apyori
Installing collected packages: apyori
Successfully installed apyori-1.1.2
(base) STM285560:~ caozeh$
```

Now, we can import the dataset ("Prac#10-GroceryStore.csv") and overview of how it is structured.

```
df = pd.read_csv("Prac#10-GroceryStore.csv", names = ['products'], sep = ',')
df.head()

products

MILK,BREAD,BISCUIT

BREAD,MILK,BISCUIT,CORNFLAKES

BREAD,TEA,BOURNVITA

JAM,MAGGI,BREAD,MILK

MAGGI,TEA,BISCUIT
```

Note: you need to download the dataset ("Prac#10-GroceryStore.csv") from Week 12 section on Learnonline course site and upload it to the Jupyter Notebook home folder.

Next, let us transform the dataset, split the items, and create a list.

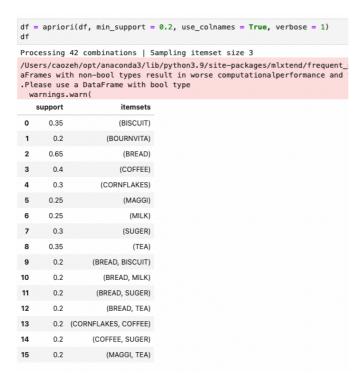
```
df.shape
data = list(df["products"].apply(lambda x:x.split(",")))
data

[['MILK', 'BREAD', 'BISCUIT'],
  ['BREAD', 'MILK', 'BISCUIT', 'CORNFLAKES'],
  ['BREAD', 'TEA', 'BOURNVITA'],
  ['JAM', 'MAGGI', 'BREAD', 'MILK'],
  ['MAGGI', 'TEA', 'BISCUIT'],
  ['MAGGI', 'TEA', 'BURNVITA'],
  ['MAGGI', 'BREAD', 'TEA', 'BISCUIT'],
  ['JAM', 'MAGGI', 'BREAD', 'TEA'],
  ['BREAD', 'MILK'],
  ['COFFEE', 'COCK', 'BISCUIT', 'CORNFLAKES'],
  ['COFFEE', 'COCK', 'BISCUIT', 'CORNFLAKES'],
  ['COFFEE', 'SUGER', 'BOURNVITA'],
  ['BREAD', 'SUGER', 'BISCUIT'],
  ['BREAD', 'SUGER', 'BOURNVITA'],
  ['BREAD', 'SUGER', 'BOURNVITA'],
  ['BREAD', 'COFFEE', 'SUGER'],
  ['BREAD', 'COFFEE', 'SUGER'],
  ['BREAD', 'COFFEE', 'SUGER'],
  ['BREAD', 'COFFEE', 'SUGER'],
  ['TEA', 'MILK', 'COFFEE', 'CORNFLAKES']]
```

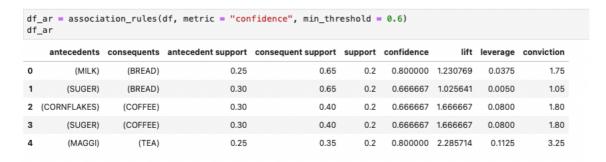
After reviewing the list items, we need to import **TransactionEncoder** that can convert the list to a One-Hot Encoded Boolean list, which is convenient for Apriori's algorithm to transform as True or False (0) values – see the next page.

<pre>from mlxtend.preprocessing import TransactionEncoder a = TransactionEncoder() a_data = a.fit(data).transform(data) df = pd.DataFrame(a_data,columns=a.columns_) df = df.replace(False,0) df</pre>											
	BISCUIT	BOURNVITA	BREAD	соск	COFFEE	CORNFLAKES	JAM	MAGGI	MILK	SUGER	TEA
0	True	0	True	0	0	0	0	0	True	0	0
1	True	0	True	0	0	True	0	0	True	0	0
2	0	True	True	0	0	0	0	0	0	0	True
3	0	0	True	0	0	0	True	True	True	0	0
4	True	0	0	0	0	0	0	True	0	0	True
5	0	True	True	0	0	0	0	0	0	0	True
6	0	0	0	0	0	True	0	True	0	0	True
7	True	0	True	0	0	0	0	True	0	0	True
8	0	0	True	0	0	0	True	True	0	0	True
9	0	0	True	0	0	0	0	0	True	0	0
10	True	0	0	True	True	True	0	0	0	0	0
11	True	0	0	True	True	True	0	0	0	0	0
12	0	True	0	0	True	0	0	0	0	True	0
13	0	0	True	True	True	0	0	0	0	0	0
14	True	0	True	0	0	0	0	0	0	True	0
15	0	0	0	0	True	True	0	0	0	True	0
16	0	True	True	0	0	0	0	0	0	True	0
17	0	0	True	0	True	0	0	0	0	True	0
18	0	0	True	0	True	0	0	0	0	True	0
19	0	0	0	0	True	True	0	0	True	0	True

The next step is to create the **Apriori Model**. We can change all the parameters in the Apriori Model in the mlxtend package (Note: there may exist some warning messages in Python due to different version upgrading issues). I will try to use minimum support parameters for this modelling - For instance, we can set a *min_support* value with a threshold value of 20% and print them on the screen as well.



The final step is that we can choose the 60% minimum confidence value. In other words, when product X is purchased, we can say that the purchase of product Y is 60% or more.



Save the <u>.html</u> document as your submission by following the instruction given on page 1 of the practical document.