

2020

Lab 8: Association Rules



ANACONDA[®]

Assoc. Prof. Dr. Mohammed Al-Sarem
Taibah University, Information System

Department

3/19/2020

Lab 8: Association Rules

Lab Objectives:

The goal of this lab is to explain the role of association rules in mining hidden relationships in data. After completion of this lab, you will be able to:

- Generate association rules using Apriori algorithm
- Evaluate the goodness of the findings rules

Methodology

First you have to install the required library that allow you to run Apriori algorithm. Then, as usual, you have to load and read dataset. After executing the algorithm, a set of generated rules are found. Your task is to analyze these rules and reports only the most interesting.

In class task:

At the end of this lab, the student will be able to:

- Explain how the Apriori algorithm works.
- Extract rules from the generated set.
- Evaluate the findings using Support and Confidence.

home task:

Complete your **Course Project** (See Home task in lab 5).

References:

- Agrawal, Rakesh, and Ramakrishnan Srikant. "Fast algorithms for mining association rules." Proc. 20th int. conf. very large data bases, VLDB. Vol. 1215. 1994.
- https://en.wikipedia.org/wiki/Apriori_algorithm

Lab 8: Association Rules

1. Frequent Itemsets via Apriori Algorithm

Apriori function to extract frequent itemsets for association rule mining. To work with apriori algorithm, first you have to install, in this time, `mlxtend` library as follows:

- add a new jupyter cell
- write down the following code: `pip install mlxtend`

Now, you can find the required algorithm by import it using the following:

```
from mlxtend.frequent_patterns import apriori
```

Generating Frequent Itemsets

Suppose we have the following transaction data as depicted in the Table below.

Table 1: Transaction Data

Tid	Items bought
1	Drink, Nuts, Diaper
2	Drink, Coffee, Diaper
3	Drink, Diaper, Eggs
4	Nuts, Eggs, Milk
5	Nuts, Coffee, Diaper, Eggs, Milk

Exercise 1.1: The apriori function expects data in a one-hot encoded pandas DataFrame. Use the transaction data that is presented in the Table 1 and convert the them into 2 dimensional array. Display the output here!

```
In [4]: dataset=[['Drink','Nuts','Diaper'],
                ['Drink','Coffee','Diaper'],
                ['Drink','Diaper','Eggs'],
                ['Nuts','Eggs','Milk'],
                ['Nuts','Coffee','Diaper','Eggs','Milk']]

In [5]: dataset
Out[5]: [['Drink', 'Nuts', 'Diaper'],
         ['Drink', 'Coffee', 'Diaper'],
         ['Drink', 'Diaper', 'Eggs'],
         ['Nuts', 'Eggs', 'Milk'],
         ['Nuts', 'Coffee', 'Diaper', 'Eggs', 'Milk']]

from mlxtend.preprocessing import TransactionEncoder

TranEncode = TransactionEncoder()
te_ary = TranEncode.fit(dataset).transform(dataset)
df = pd.DataFrame(te_ary, columns=TranEncode.columns_)
```

Lab 8: Association Rules

Exercise 1.2: Display the dataset below! What do you see?

```
In [6]: import pandas as pd
        from mlxtend.preprocessing import TransactionEncoder

        TranEncod = TransactionEncoder()
        te_ary = TranEncod.fit(dataset).transform(dataset)
        df = pd.DataFrame(te_ary, columns=TranEncod.columns_)
```

In [7]: df

Out[7]:

	Coffee	Diaper	Drink	Eggs	Milk	Nuts
0	False	True	True	False	False	True
1	True	True	True	False	False	False
2	False	True	True	True	False	False
3	False	False	False	True	True	True
4	True	True	False	True	True	True

*what i see is
it show us if
the item exist
in the dataset
or not by True
and False*

Remember that to work with association rules, the algorithm required as input beside the dataset, the minimum support value ϵ :

```
Apriori( $T, \epsilon$ )
 $L_1 \leftarrow \{\text{large 1-itemsets}\}$ 
 $k \leftarrow 2$ 
while  $L_{k-1} \neq \emptyset$ 
     $C_k \leftarrow \{c = a \cup \{b\} \mid a \in L_{k-1} \wedge b \notin a, \{s \subseteq c \mid |s| = k-1\} \subseteq L_{k-1}\}$ 
    for transactions  $t \in T$ 
         $D_t \leftarrow \{c \in C_k \mid c \subseteq t\}$ 
        for candidates  $c \in D_t$ 
             $\text{count}[c] \leftarrow \text{count}[c] + 1$ 
     $L_k \leftarrow \{c \in C_k \mid \text{count}[c] \geq \epsilon\}$ 
     $k \leftarrow k + 1$ 
return  $\bigcup_k L_k$ 
```

So, let us return the items and itemsets with at least 60% support.

```
apriori(df, min_support=0.6)
```

By default, `apriori` returns the column indices of the items, which may be useful in downstream operations such as association rule mining. For better readability, we can set `use_colnames=True` to convert these integer values into the respective item names:

```
: apriori(df, min_support=0.6, use_colnames=True)
```

	support	itemsets
0	0.8	(Diaper)
1	0.6	(Drink)
2	0.6	(Eggs)
3	0.6	(Nuts)
4	0.6	(Diaper, Drink)

Lab 8: Association Rules

Exercise 1.3: Find the possible rules and calculate the confidence of each rule you find?

```
from mxtext.frequent_patterns import association_rules
```

```
association_rules(frequent_itemsets, metric="confidence", min_threshold=0.7)
```

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	leverage	conviction
0	(Drink)	(Diaper)	0.6	0.8	0.6	1.00	1.25	0.12	inf
1	(Diaper)	(Drink)	0.8	0.6	0.6	0.75	1.25	0.12	1.6

2. Selecting and Filtering Results

Before moving ahead, let us decrease the minimum support to 30%. This leads to increase the number of returned items and itemsets. Let us now filter out these rules. The advantage of working with pandas **DataFrames** is that we can use its convenient features to filter the results. For instance, let's assume we are only interested in **itemsets** of length 2 that have a support of at least 50% percent. First, we create the frequent itemsets via **apriori** and add a new column that stores the length of each itemset:

```
frequent_itemsets = apriori(df, min_support=0.5, use_colnames=True)
frequent_itemsets['length'] = frequent_itemsets['itemsets'].apply(lambda x: len(x))
frequent_itemsets
```

	support	itemsets	length
0	0.8	(Diaper)	1
1	0.6	(Drink)	1
2	0.6	(Eggs)	1
3	0.6	(Nuts)	1
4	0.6	(Diaper, Drink)	2

Then, you can select the results that satisfy our desired criteria as follows:

```
frequent_itemsets[ (frequent_itemsets['length'] == 2) &
                   (frequent_itemsets['support'] >= 0.5) ]
```

	support	itemsets	length
4	0.6	(Diaper, Drink)	2

Similarly, using the Pandas API, we can select entries based on the "itemsets" column:

```
frequent_itemsets[ frequent_itemsets['itemsets'] == {'Diaper', 'Drink'} ]
```

Good luck