VIETNAMESE - GERMAN UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE AND ENGINEER

# Lab 5 Report

*Student:*    Le Thanh Hai - 10421016

*Student:*    Trinh The Hao - 10421017

*Student:*    Vu Ngoc Quang - 10421103

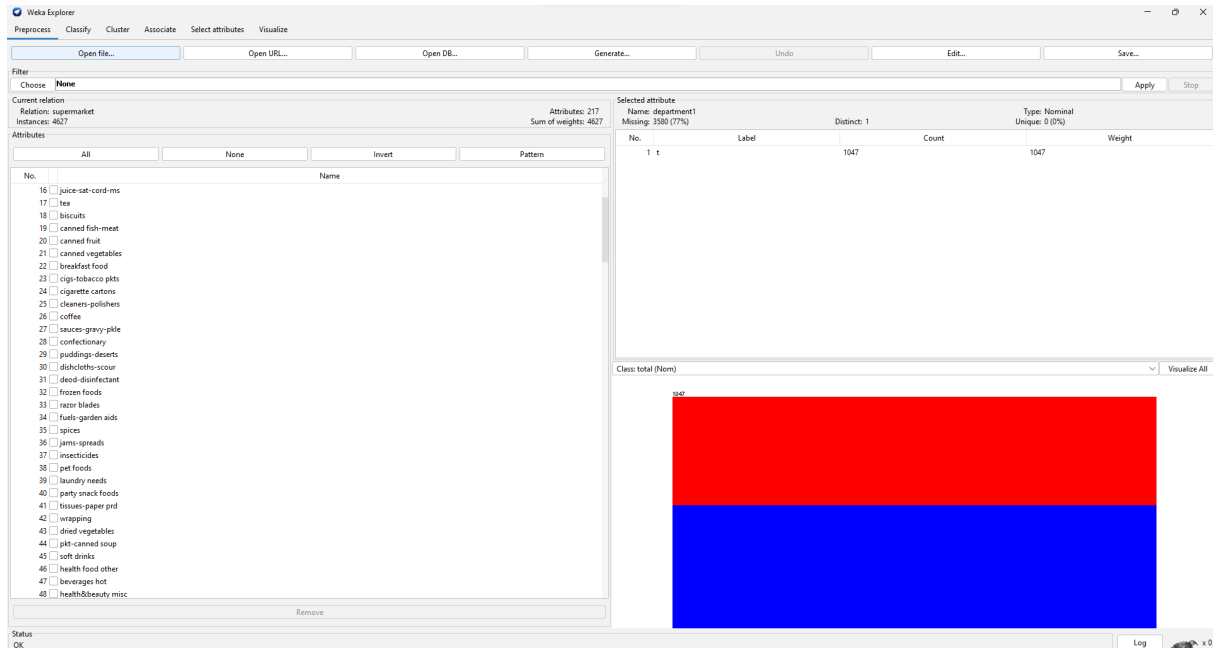07 / 2022

# 1 Exercise 2.1,2,3,4,5,6
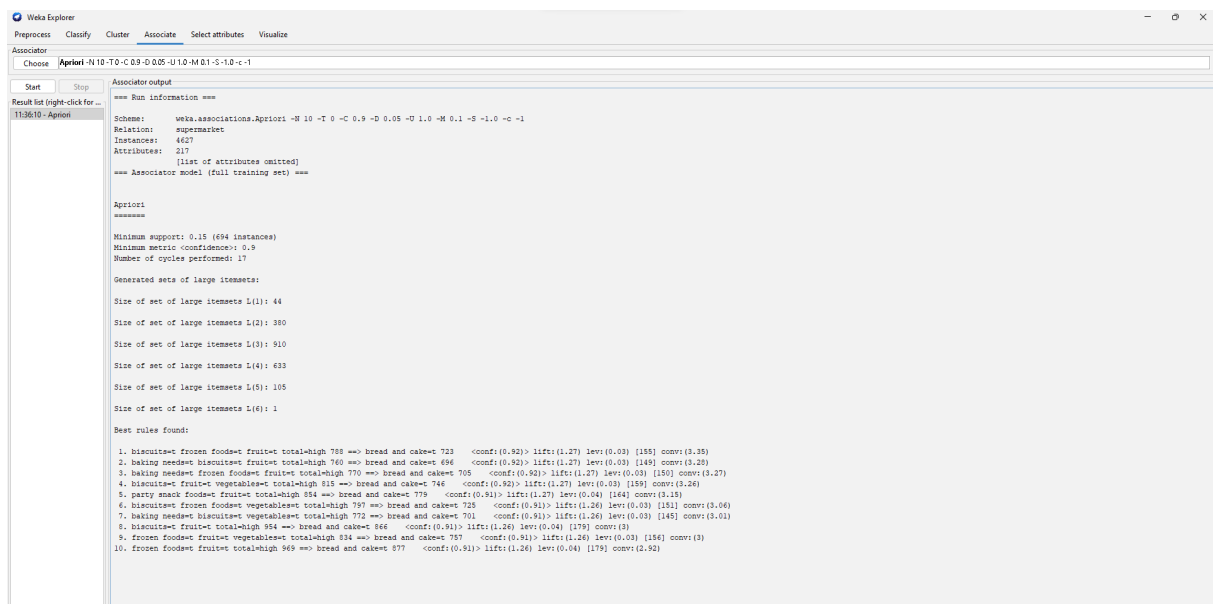


Figure 1: 2.1 load



Figure 2: 2.2 describe

**Brief Explanation:**

To load data ,In the WEKA explorer, open the Preprocess tab, click on the Open file ... button and select supermarket.arff database from the installation folder.Click on the Associate
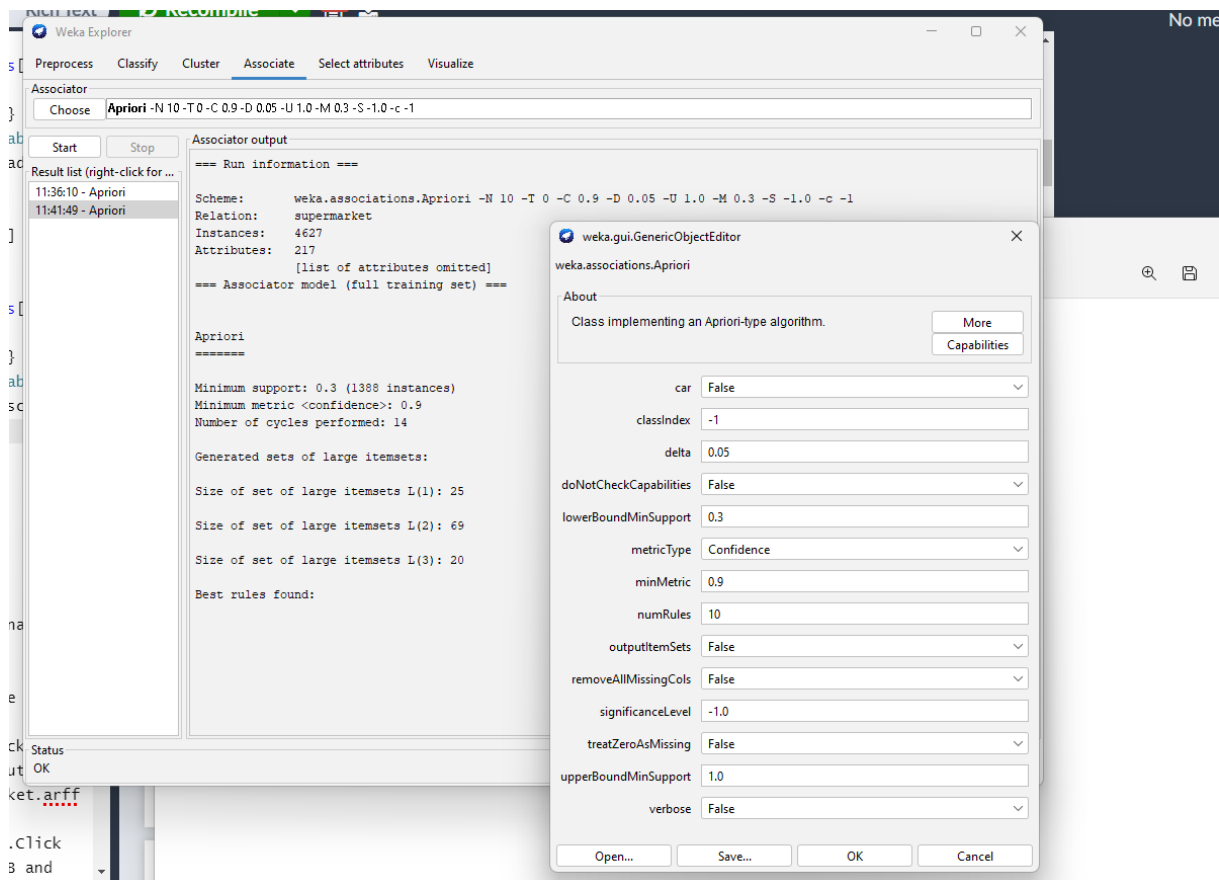
Figure 3: 2.3 2.4 2.5 2.6

TAB and click on the Choose button. Select the Apriori association .

To describe the dataset, you need to load it in Weka and examine its attributes and instances. In Weka, you can open the dataset file (supermarket.arff) using the "Explorer" interface. The number of instances refers to the number of rows or examples in the dataset, while the number of items refers to the number of attributes or columns in the dataset.

Run Apriori algorithm with minsup=0.3 and minconf=0.9: In Weka, after loading the dataset, you can navigate to the "Associate" tab in the "Explorer" interface. Select the "Apriori" algorithm and set the minimum support (minsup) to 0.3 and the minimum confidence (minconf) to 0.9. Run the algorithm, and Weka will generate association rules based on these parameters.

Report top ten rules: Once the Apriori algorithm finishes running, Weka will display the generated association rules. You can sort them based on various metrics such as confidence, support, or lift. To report the top ten rules, you can sort the rules based on the desired metric and select the top ten rules with the highest values.

Run Apriori with different minsup and minconf: To run the Apriori algorithm with different minimum support and minimum confidence values, you can repeat the process described in step 2.3. Simply change the minsup and minconf parameters to the desired values, and Weka will generate new association rules based on these updated parameters.

Report top ten rules: Similarly, after running the updated Apriori algorithm, you can sort and report the top ten rules based on the desired metric as explained in step 2.4.

Please note that the specific steps and options may vary slightly depending on the version of Weka you are using. It's always a good idea to consult the Weka documentation or user guides for detailed instructions on how to perform these tasks within the software.

## 2   Exercise 2.7 & 2.8 & 2.9 & 2.10

```python
from scipy.io import arff
import pandas as pd
from scipy.io import arff
import apyori
from apyori import apriori

# Load the data set.
data = arff.loadarff("supermarket.arff")
df = pd.DataFrame(data[0])
```

```
10
11  # drop the total column
12  df = df.drop(['total'],axis=1)
13
14  # convert the data to a list of attributes
15  attributes = list(df.columns)
16
17  # transform the data into strings and into a numpy matrix for easier ↩
        handling
18  sales = df.to_numpy().astype(str)
19
20  transactions = [ [] for _ in range(df.shape[0]) ]  #Empty lists to be ↩
        filled with the specific products purchased.
21
22  #We iterate over every element in the sales matrix and, if a product ↩
        belongs to a transaction, we save it.
23  for i in range(df.shape[0]):
24      for j in range(df.shape[1]):
25          if sales[i,j] == "t":
26              transactions[i].append(attributes[j])
27
28
29  mins = 0.1 # minimum support
30  minc = 0.9 # minimum confidence
31
32  # We use the apriori algorithm to find the association rules.
33  results = list(apriori(transactions,min_support=mins,min_confidence=↩
        minc))
34
35  for i in range(0,10):
36      print(list(results[i][2][0][0]),'->',list(results[i][2][0][1]))
37
38  print('There are '+str(len(results))+' rules with a minimum support of↩
        '+str(mins)+
39      ', a minimum confidence of ' +str(minc))
```

**Brief Explanation:** The code loads a supermarket dataset from an ARFF file and converts it into a pandas DataFrame. It removes the "total" column from the DataFrame. The data is transformed into a list of transactions, where each transaction represents a list of products purchased. The code defines minimum support (mins) and minimum confidence (minc) thresholds. The Apriori algorithm is applied to the transactions list to find association rules that meet the minimum support and minimum confidence criteria. The code prints the first 10 association rules found, showing the antecedent (products purchased together) and consequent (product purchased after the antecedent). Finally, the code prints the total number of association rules found with the

specified minimum support and minimum confidence thresholds.

```
['biscuits', 'milk-cream', 'canned fruit'] -> ['bread and cake']
['small goods', 'biscuits', 'frozen foods'] -> ['bread and cake']
['small goods', 'biscuits', 'milk-cream'] -> ['bread and cake']
['tissues-paper prd', 'milk-cream', 'canned fruit'] -> ['bread and cake']
['baking needs', 'beef', 'milk-cream', 'biscuits'] -> ['bread and cake']
['tissues-paper prd', 'baking needs', 'beef', 'milk-cream'] -> ['bread and cake']
['laundry needs', 'baking needs', 'milk-cream', 'fruit'] -> ['bread and cake']
['margarine', 'baking needs', 'fruit', 'pet foods'] -> ['bread and cake']
['laundry needs', 'baking needs', 'tissues-paper prd', 'milk-cream'] -> ['bread and cake']
['margarine', 'baking needs', 'milk-cream', 'pet foods'] -> ['bread and cake']
There are 100 rules with a minimum support of 0.1, a minimum confidence of 0.9
```

# 3    Exercise 2.11 & 2.12

```python
from scipy.io import arff
import pandas as pd
from scipy.io import arff
import apyori
from apyori import apriori


# Load the data set.
data = arff.loadarff("supermarket.arff")
df = pd.DataFrame(data[0])


# drop the total column
df = df.drop(['total'],axis=1)


# convert the data to a list of attributes
attributes = list(df.columns)


# transform the data into strings and into a numpy matrix for easier ↩
    handling
sales = df.to_numpy().astype(str)

transactions = [ [] for _ in range(df.shape[0]) ]  #Empty lists to be ↩
    filled with the specific products purchased.

#We iterate over every element in the sales matrix and, if a product ↩
    belongs to a transaction, we save it.
for i in range(df.shape[0]):
    for j in range(df.shape[1]):
        if sales[i,j] == "t":
            transactions[i].append(attributes[j])

```

```python
28
29  mins = [0.1,0.2,0.3] # minimum support
30  minc = 0.9 # minimum confidence
31
32  for m in mins:
33      # We use the apriori algorithm to find the association rules.
34      results = list(apriori(transactions,min_support=m,min_confidence=↵
            minc))
35      if len(results) > 10:
36          for i in range(0,10):
37              print(list(results[i][2][0][0]),'->',list(results[i↵
                    ][2][0][1]))
38      else:
39          for i in range(0,len(results)):
40              print(list(results[i][2][0][0]),'->',list(results[i↵
                    ][2][0][1]))
41
42      print('There are '+str(len(results))+' rules with a minimum ↵
            support of '+str(m)+
43          ', a minimum confidence of ' +str(minc))
44      print('')
45
46  mins = 0.1 # minimum support
47  minc = [0.9,0.8,0.7] # minimum confidence
48
49  for m in minc:
50      # We use the apriori algorithm to find the association rules.
51      results = list(apriori(transactions,min_support=mins,↵
            min_confidence=m))
52      if len(results) > 10:
53          for i in range(0,10):
54              print(list(results[i][2][0][0]),'->',list(results[i↵
                    ][2][0][1]))
55      else:
56          for i in range(0,len(results)):
57              print(list(results[i][2][0][0]),'->',list(results[i↵
                    ][2][0][1]))
58
59      print('There are '+str(len(results))+' rules with a minimum ↵
            support of '+str(mins)+
60          ', a minimum confidence of ' +str(m))
61      print('')
```

```
['canned fruit', 'milk-cream', 'biscuits'] -> ['bread and cake']
['small goods', 'frozen foods', 'biscuits'] -> ['bread and cake']
['small goods', 'milk-cream', 'biscuits'] -> ['bread and cake']
['tissues-paper prd', 'canned fruit', 'milk-cream'] -> ['bread and cake']
['baking needs', 'milk-cream', 'beef', 'biscuits'] -> ['bread and cake']
['tissues-paper prd', 'baking needs', 'beef', 'milk-cream'] -> ['bread and cake']
['fruit', 'baking needs', 'laundry needs', 'milk-cream'] -> ['bread and cake']
['fruit', 'baking needs', 'margarine', 'pet foods'] -> ['bread and cake']
['tissues-paper prd', 'baking needs', 'laundry needs', 'milk-cream'] -> ['bread and cake']
['baking needs', 'milk-cream', 'margarine', 'pet foods'] -> ['bread and cake']
There are 100 rules with a minimum support of 0.1, a minimum confidence of 0.9

There are 0 rules with a minimum support of 0.2, a minimum confidence of 0.9

There are 0 rules with a minimum support of 0.3, a minimum confidence of 0.9

['canned fruit', 'milk-cream', 'biscuits'] -> ['bread and cake']
['small goods', 'frozen foods', 'biscuits'] -> ['bread and cake']
['small goods', 'milk-cream', 'biscuits'] -> ['bread and cake']
['tissues-paper prd', 'canned fruit', 'milk-cream'] -> ['bread and cake']
['baking needs', 'milk-cream', 'beef', 'biscuits'] -> ['bread and cake']
['tissues-paper prd', 'baking needs', 'beef', 'milk-cream'] -> ['bread and cake']
['fruit', 'baking needs', 'laundry needs', 'milk-cream'] -> ['bread and cake']
['fruit', 'baking needs', 'margarine', 'pet foods'] -> ['bread and cake']
['tissues-paper prd', 'baking needs', 'laundry needs', 'milk-cream'] -> ['bread and cake']
['baking needs', 'milk-cream', 'margarine', 'pet foods'] -> ['bread and cake']
There are 100 rules with a minimum support of 0.1, a minimum confidence of 0.9

['canned fish-meat'] -> ['bread and cake']
['canned fruit'] -> ['bread and cake']
['jams-spreads'] -> ['bread and cake']
['margarine'] -> ['bread and cake']
['potatoes'] -> ['bread and cake']
['puddings-deserts'] -> ['bread and cake']
['small goods'] -> ['bread and cake']
['potatoes'] -> ['vegetables']
['baking needs', 'beef'] -> ['bread and cake']
['baking needs', 'biscuits'] -> ['bread and cake']
There are 3672 rules with a minimum support of 0.1, a minimum confidence of 0.8

[] -> ['bread and cake']
['baby needs'] -> ['bread and cake']
['baking needs'] -> ['bread and cake']
['breakfast food'] -> ['baking needs']
['canned fish-meat'] -> ['baking needs']
['canned fruit'] -> ['baking needs']
['canned vegetables'] -> ['baking needs']
['cleaners-polishers'] -> ['baking needs']
['coffee'] -> ['baking needs']
['dental needs'] -> ['baking needs']
There are 6873 rules with a minimum support of 0.1, a minimum confidence of 0.7
```
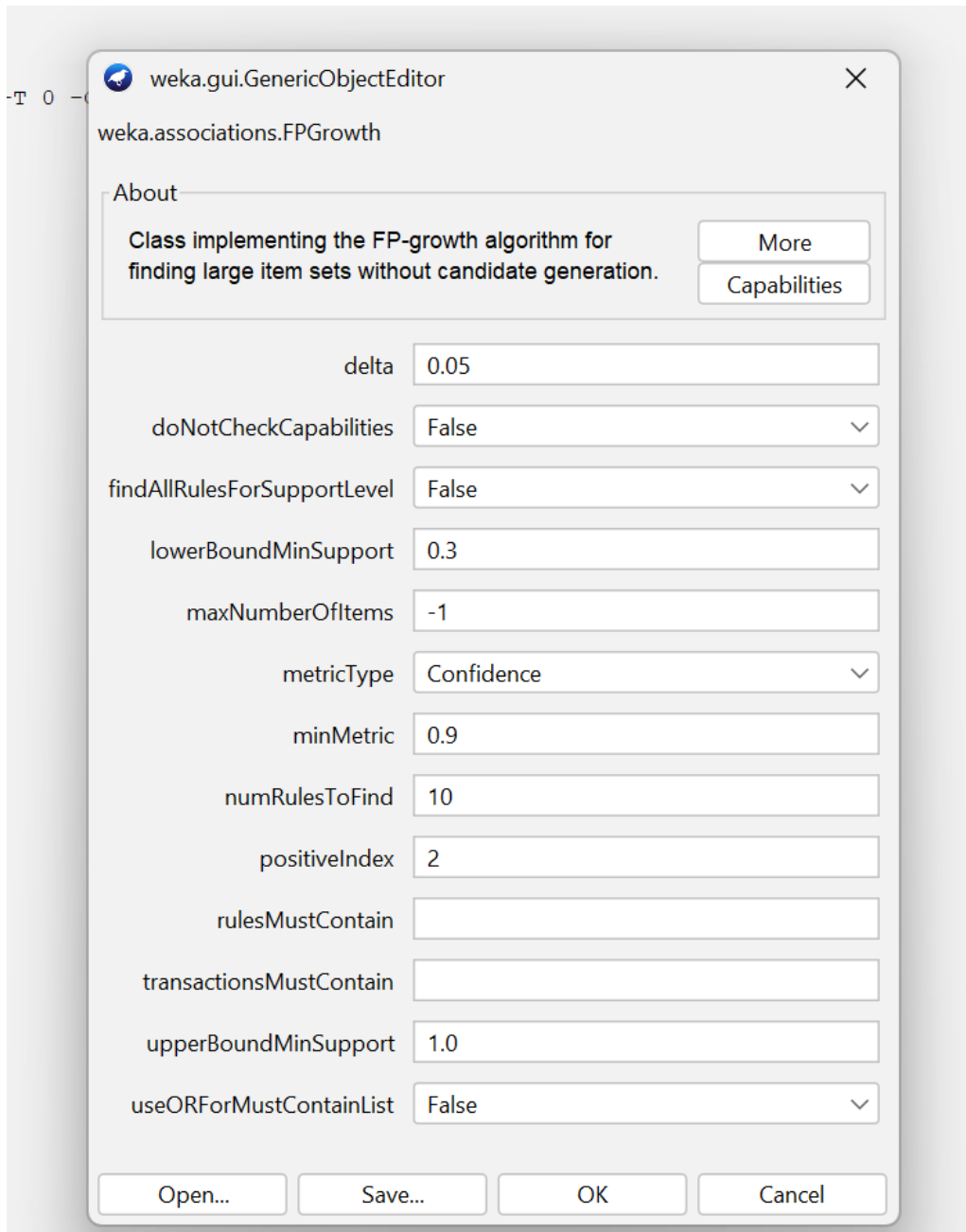
## 4   Exercise 2.13 & 2.14

```
-T 0 -(
```

**weka.gui.GenericObjectEditor**                                        ✕

weka.associations.FPGrowth

**About**

Class implementing the FP-growth algorithm for     [ More ]
finding large item sets without candidate generation.   [ Capabilities ]

| | |
|---|---|
| delta | 0.05 |
| doNotCheckCapabilities | False |
| findAllRulesForSupportLevel | False |
| lowerBoundMinSupport | 0.3 |
| maxNumberOfItems | -1 |
| metricType | Confidence |
| minMetric | 0.9 |
| numRulesToFind | 10 |
| positiveIndex | 2 |
| rulesMustContain | |
| transactionsMustContain | |
| upperBoundMinSupport | 1.0 |
| useORForMustContainList | False |

[ Open... ]   [ Save... ]   [ OK ]   [ Cancel ]

```
=== Run information ===

Scheme:       weka.associations.FPGrowth -P 2 -I -1 -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.3
Relation:     supermarket
Instances:    4627
Attributes:   217
              [list of attributes omitted]
=== Associator model (full training set) ===

No rules found!
```

**Description:**

We tried to set the minsup value to 0.3 and the mincof value to 0.9 and the output we got was there were no rules found. We will try to figure out what happened by setting the difference values in the next part.

# 5    Exercise 2.15 & 2.16

```
=== Run information ===

Scheme:       weka.associations.FPGrowth -P 2 -I -1 -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1
Relation:     supermarket
Instances:    4627
Attributes:   217
              [list of attributes omitted]
=== Associator model (full training set) ===


FPGrowth found 16 rules (displaying top 10)

 1. [fruit=t, frozen foods=t, biscuits=t, total=high]: 788 ==> [bread and cake=t]: 723    <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.35)
 2. [fruit=t, baking needs=t, biscuits=t, total=high]: 760 ==> [bread and cake=t]: 696    <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.28)
 3. [fruit=t, baking needs=t, frozen foods=t, total=high]: 770 ==> [bread and cake=t]: 705    <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.27)
 4. [fruit=t, vegetables=t, biscuits=t, total=high]: 815 ==> [bread and cake=t]: 746    <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.26)
 5. [fruit=t, party snack foods=t, total=high]: 854 ==> [bread and cake=t]: 779    <conf:(0.91)> lift:(1.27) lev:(0.04) conv:(3.15)
 6. [vegetables=t, frozen foods=t, biscuits=t, total=high]: 797 ==> [bread and cake=t]: 725    <conf:(0.91)> lift:(1.26) lev:(0.03) conv:(3.06)
 7. [vegetables=t, baking needs=t, biscuits=t, total=high]: 772 ==> [bread and cake=t]: 701    <conf:(0.91)> lift:(1.26) lev:(0.03) conv:(3.01)
 8. [fruit=t, biscuits=t, total=high]: 954 ==> [bread and cake=t]: 866    <conf:(0.91)> lift:(1.26) lev:(0.04) conv:(3)
 9. [fruit=t, vegetables=t, frozen foods=t, total=high]: 834 ==> [bread and cake=t]: 757    <conf:(0.91)> lift:(1.26) lev:(0.03) conv:(3)
10. [fruit=t, frozen foods=t, total=high]: 969 ==> [bread and cake=t]: 877    <conf:(0.91)> lift:(1.26) lev:(0.04) conv:(2.92)
```
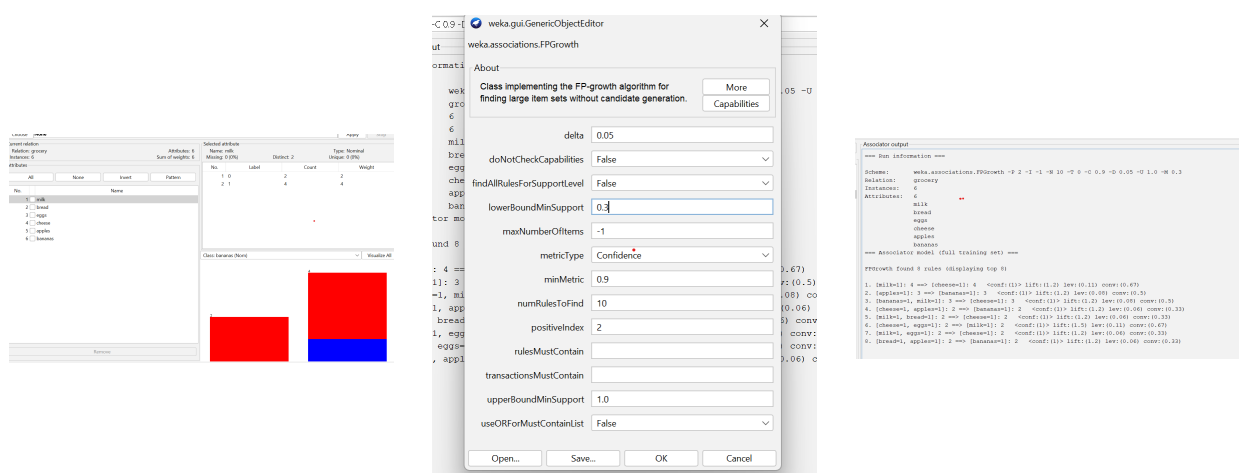
**Description:**

We tried to set a new value for the minsup to 0.1 and this time it did show out the top ten rules.

It turned out that When you set minsup to 0.3, the algorithm will only consider itemsets that appear in at least 30 percent of the transactions. If there are no itemsets that appear in at least 30 percent of the transactions, then the algorithm will not find any rules.
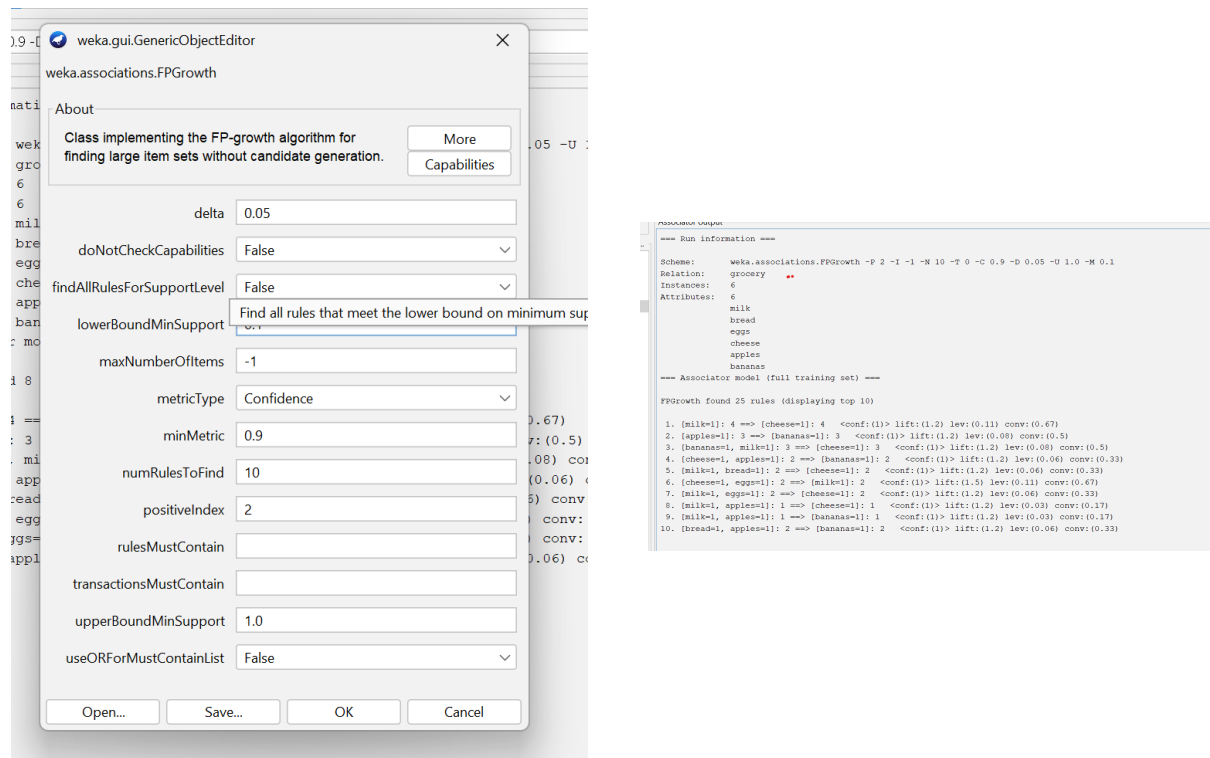
When you set minsup to 0.1, the algorithm will consider itemsets that appear in at least 10 percent of the transactions. This means that the algorithm is more likely to find rules, even if the rules are not very strong.

# 6 Exercise 2.17 & 2.18



**Description: We created our own dataset and the first one was with minsup =0.3**

**Description: We then tried it with other minsup =0.1**

# 7 Exercise 2.19 & 2.20

```
1  # Sample dataset for movie recommendations
2  dataset = {
3      'movie1': {'action': 1, 'thriller': 1, 'comedy': 0, 'romance': 0,'↩
           adventure':0},
4      'movie2': {'action': 0, 'thriller': 0, 'comedy': 1, 'romance': 1,'↩
           adventure':0},
5      'movie3': {'action': 1, 'thriller': 0, 'comedy': 0, 'romance': 0,'↩
           adventure':1},
6      'movie4': {'action': 0, 'thriller': 0, 'comedy': 1, 'romance': 0,'↩
           adventure':0},
7      'movie5': {'action': 0, 'thriller': 1, 'comedy': 0, 'romance': 0,'↩
           adventure':1},
8      'movie6': {'action': 1, 'thriller': 0, 'comedy': 1, 'romance': 0,'↩
           adventure':0},
9      'movie7': {'action': 0, 'thriller': 0, 'comedy': 0, 'romance': 1,'↩
           adventure':0},
10     'movie8': {'action': 0, 'thriller': 1, 'comedy': 0, 'romance': 0,'↩
           adventure':0},
11     'movie9': {'action': 1, 'thriller': 0, 'comedy': 0, 'romance': 0,'↩
           adventure':0},
12     'movie10': {'action': 0, 'thriller': 0, 'comedy': 1, 'romance': 1,↩
           'adventure':0}
```

```python
13  }
14
15  def content_based_recommender(dataset, user_preferences, ↩
        num_recommendations=5):
16      # Compute weighted sums
17      scores = {}
18      for movie in dataset:
19          features = dataset[movie]
20          weighted_sum = sum(features[genre] for genre in ↩
                user_preferences)
21          scores[movie] = weighted_sum
22
23      # Sort the scores in descending order
24      sorted_scores = sorted(scores.items(), key=lambda x: x[1], reverse↩
            =True)
25
26      # Generate recommendations
27      recommendations = []
28      for movie, score in sorted_scores:
29          if movie not in user_preferences:
30              recommendations.append(movie)
31          if len(recommendations) == num_recommendations:
32              break
33
34      return recommendations
35
36  # Example usage
37  user_preferences = ['action', 'adventure']
38  recommendations = content_based_recommender(dataset, user_preferences)
39  print("Recommendations:", recommendations)
```

```
Recommendations: ['movie3', 'movie1', 'movie5', 'movie6', 'movie9']
```