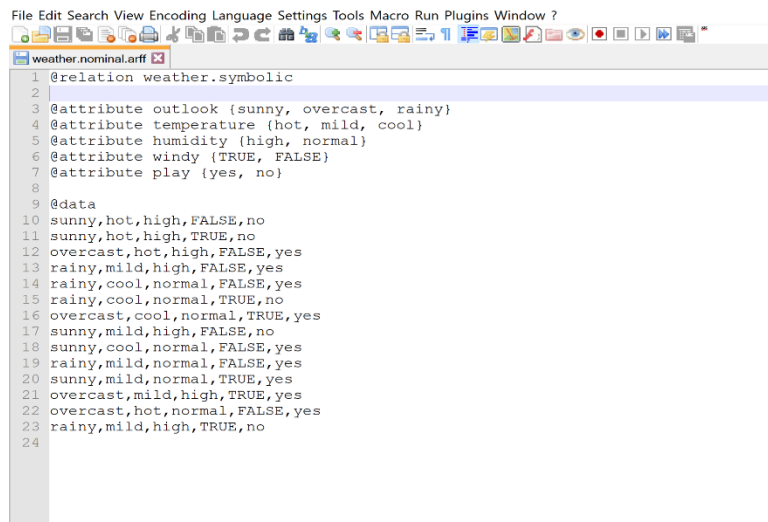# Solution to Laboratory Week 5

**1. You will need to have access to the WEKA package.**

**2. The data files for this lab can be found at /KDrive/SEH/SCSIT/Students/Courses/COSC2111/DataMining/data**

**3. Load the file arff/UCI/weather.nominal.arff.**

**(a) Inspect the file with an editor.**



Figure 1. Snapshot of view of weather.nominal data set in Notpad++(editor).

As it can be seen from figure1, the data set contains 5 categorical attributes (i.e., 5 columns) and 12 instances (i.e., 12 rows). The last attribute with the name of "play" is the target value and has two classes/labels of yes and no.

**(b) Run Apriori using the default settings of the options, but turn on the printing of item sets.**
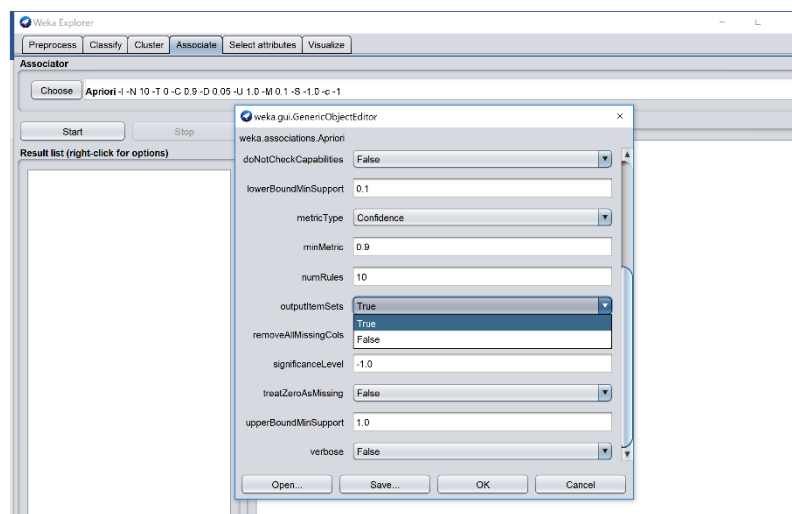


Figure 2. View of selecting Apriori algorithm from Associate tab and how to turn on the printing of item sets by selecting True on OutputItemSets menu.

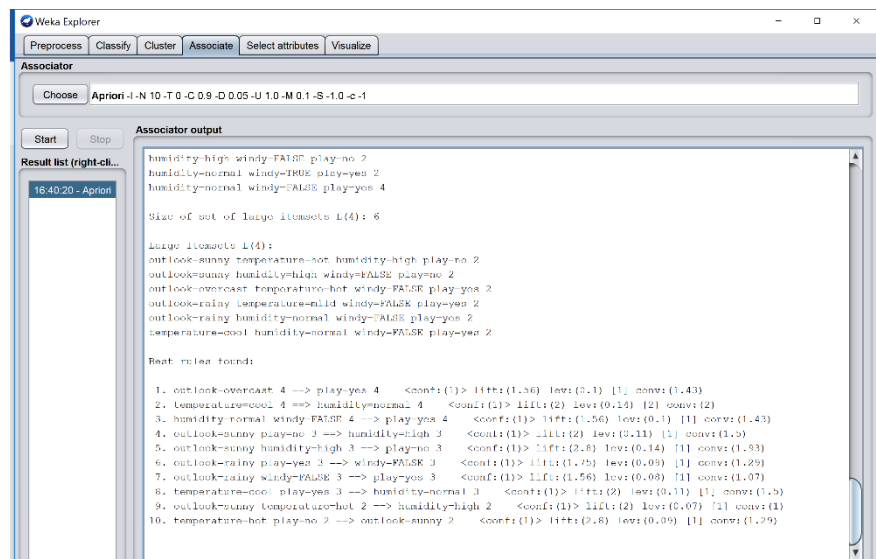**(c) What is the confidence for rule 10. How was this confidence value computed?**



Figure 3. Output from Run of Apriori on weather.nominal data set.

Below is rule 10, where its confidence is 1 as shown conf:(1).

Rule : "10. temperature=hot play=no 2 ==> outlook=sunny 2    <conf:(1)> lift:(2.8) lev:(0.09) [1] conv:(1.29)"

Support is 2 or 2/14 or 14.2%

Confidence is 2/2 = 1 or 100%

**(d) How many instances form the support of rule 8?**

3 instances.

**(e) Identify the itemset from which rule 1 was generated. What other rules could be generated from this itemset?**

This is rule 1 :  outlook=overcast 4 ==> play=yes 4.

The itemset that the rule 1 has been generated from is a set contains outlook=overcast and play= yes items (i.e., set of rule 1 = {outlook=overcast, play= yes}.

**(f) What does 'best rules' mean? What criteria are used to determine the best rules?**

The best rules are rules with the highest support and high confidence (i.e., confidence equal or close to 1).

**(g) How many rules can be generated from this data? Experiment with the numRules parameter.**

As shown in figure 4, by setting numRules to 2000, the maximum number of rules that will pull out by running Apriori on top of the data set is 336 rules. You can try with smaller or larger number than 336 to see the difference. If you set number of rules to any number greater than 336, it would only return 336 rules not any more.
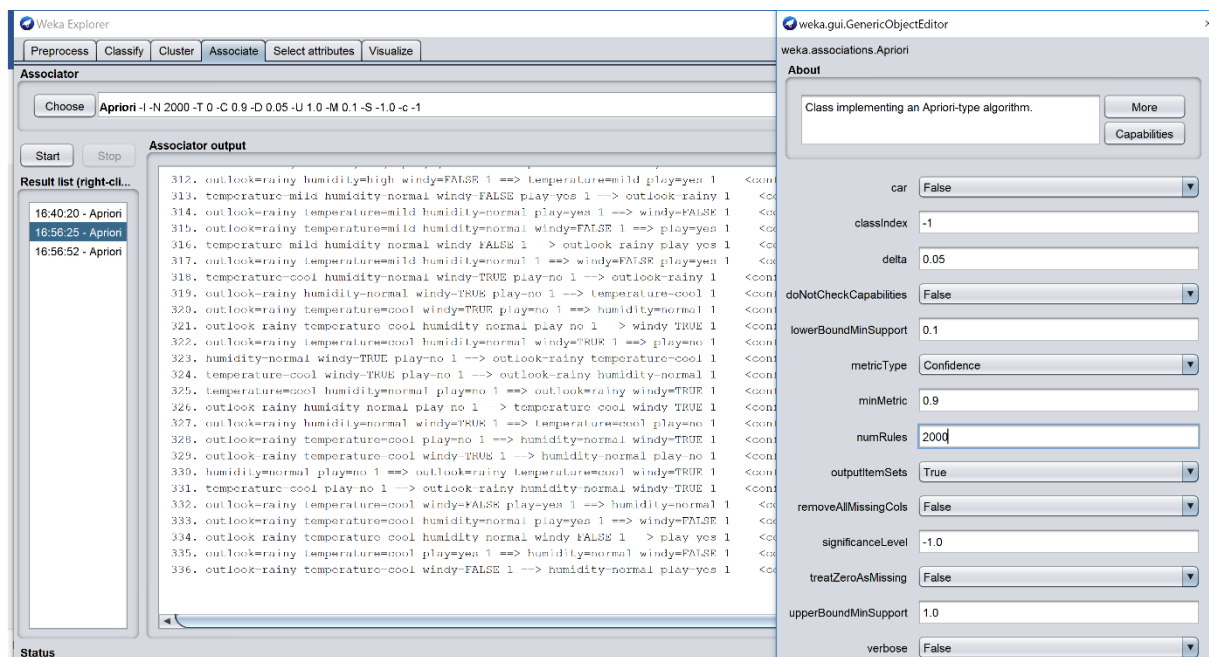
Figure 4. Maximum number of rules that pules out from running apriori on weather data set is 336 even by setting parameter numRules= 2000.

**4. Load the file arff/supermarket1-subset.arff.**

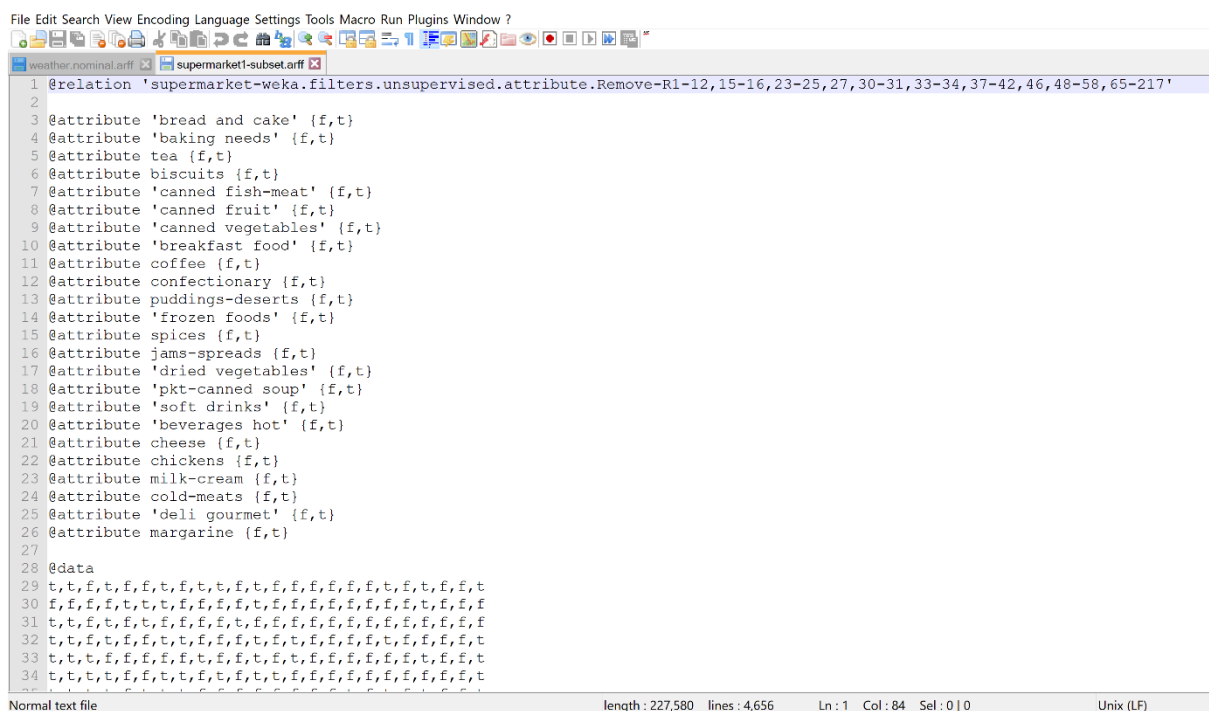**(a) View the file with an editor.**



Figure 5. View of the data set in Notepad++ editor.

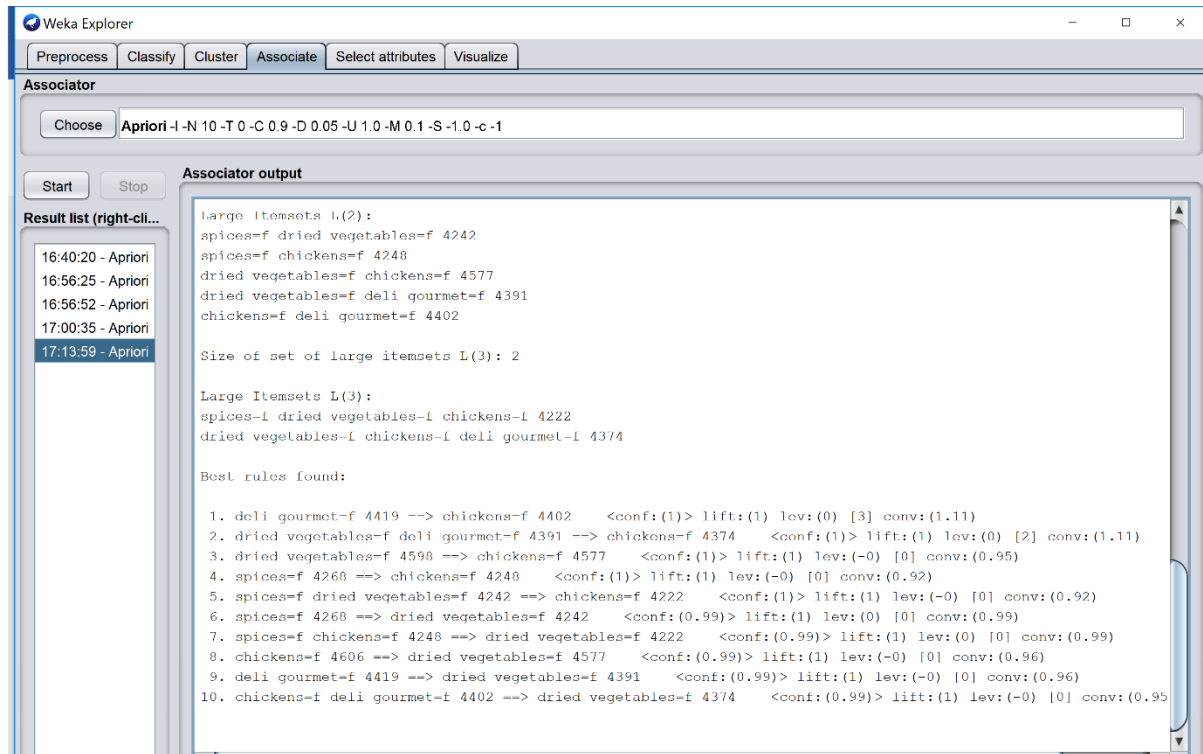**(b) Go to the Associate screen and run Apriori with the default values.**



Figure 6. Snapshot of outputs of running Apriori on supermarket-subset dataset.

**(c) Are there any rules with high accuracy?**

Many

**(d) Are there any golden nuggets in the output?**

There are many strong rules, however the relationships captured are of the form: not buying some items results in not buying other items. These relationships are not useful in practice and are unlikely to yield golden nuggets.

(Questions e, f, g are editing mistakes.)

**(h) Experiment with different metrics. How do the generated rules change with the different metrics.**
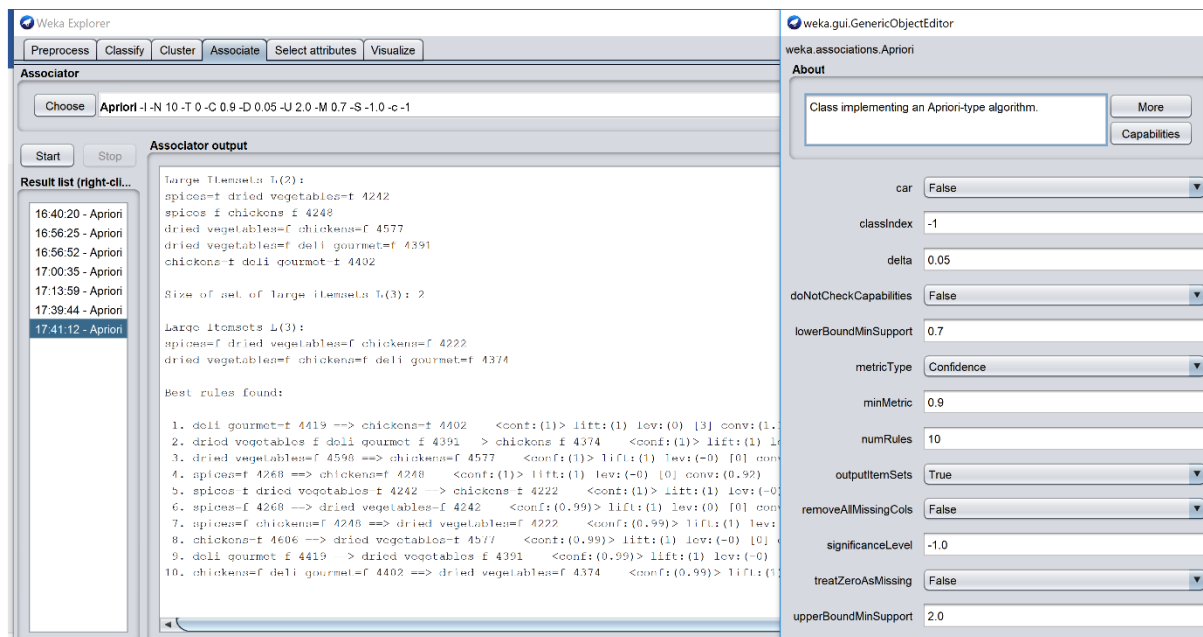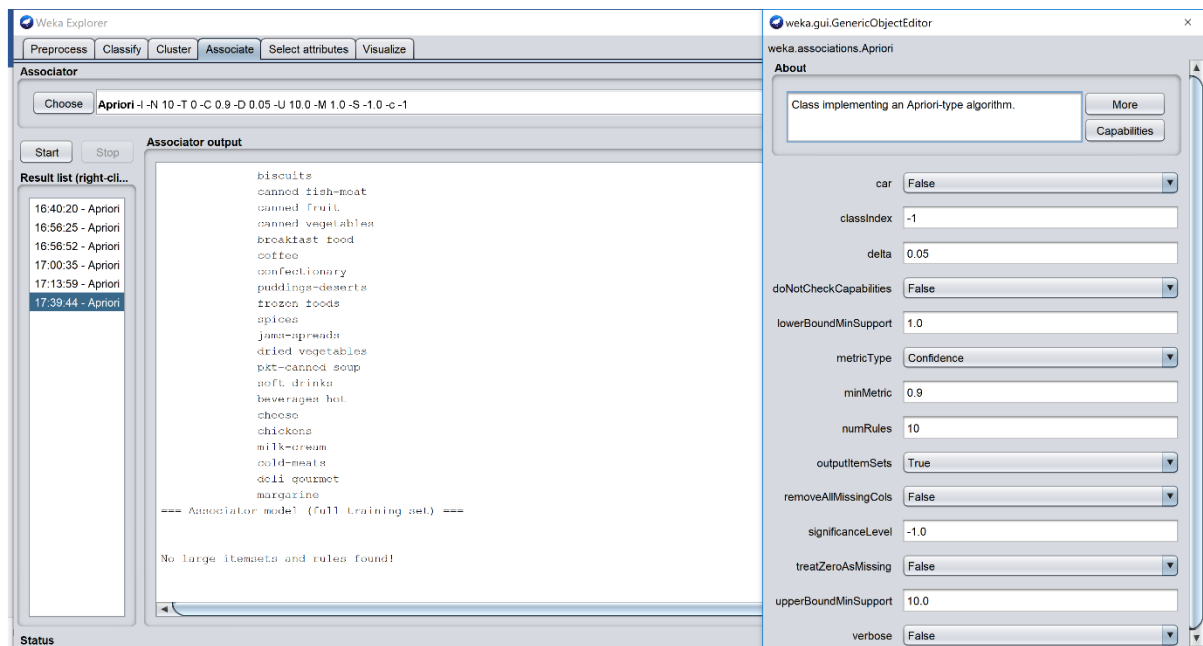
Figure 7. Effect of changing lowerBoundMinSupport and UpperBoundMinSupport.

**5. Load the file arff/supermarket2-small.arff.**
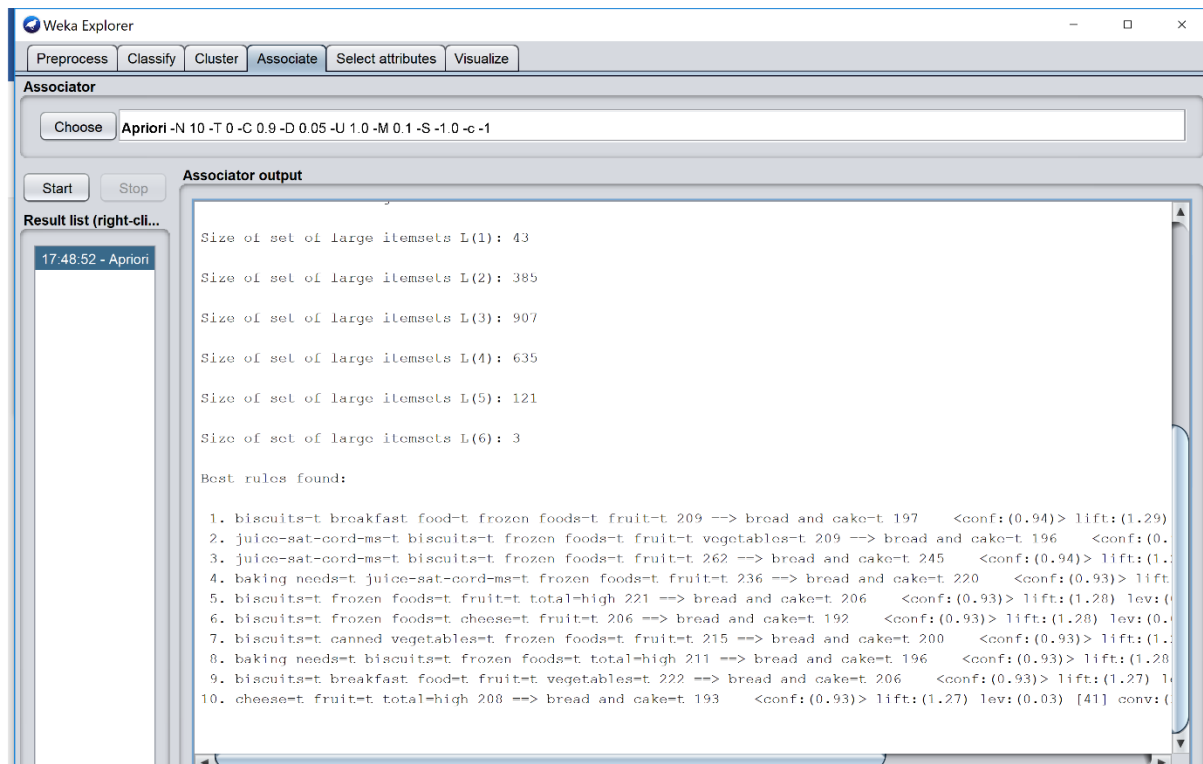
**(a) Repeat (4) above on this data.**



Figure 8. Output of running Apriori on supermarket2-small data set.s

**(b) What do you conclude about the different representations?**

In previous data set all values are represented by true or false as "t" (i.e., purchased) and "f"(i.e., not purchased) and  apriori generated rules from not purchased items which were not very useful.  In this data set, values are represented by either "t" or "?"  (i.e., missing values) and the generated rules are between purchased items.   Most of the strongest rules concern bread and cake.  I don't know enough about supermarket operations to know whether this is a golden nugget or not.

**6. Repeat (4) and (5) with FilteredAssociator and FPGrowth and compare the results with apriori.**

**7. The files supermarket1-subset.arff and supermarket2-small.arff are trimmed down versions of supermarket1.arff and supermarket2.arff in order to get quick execution times for apriori.**

**(a) Run apriori on the larger files.**

There are memory and computation time problems.  How long should we wait for the run to finish? Hours? Days? Weeks?

**(b) What problems arise and what can you do about them?**

We need to experiment with smaller numbers of instances and/or attributes to get a sense of how the algorithms scale with increasing amounts of data.
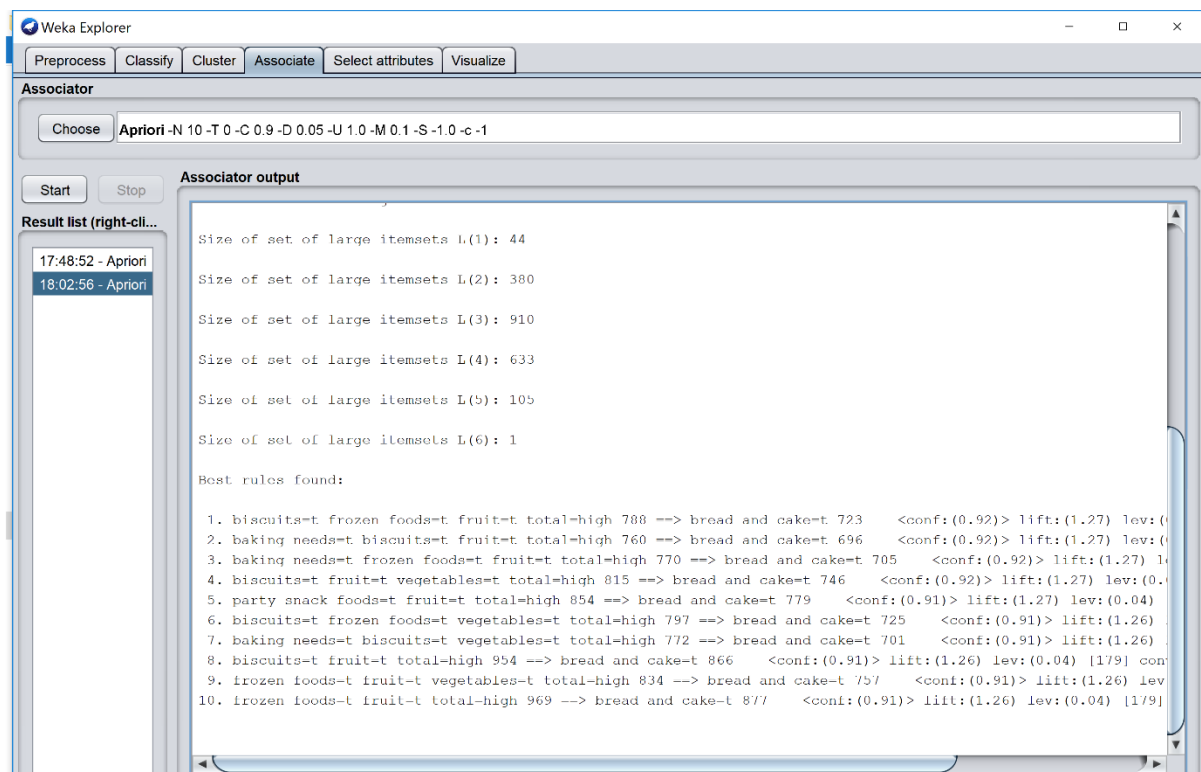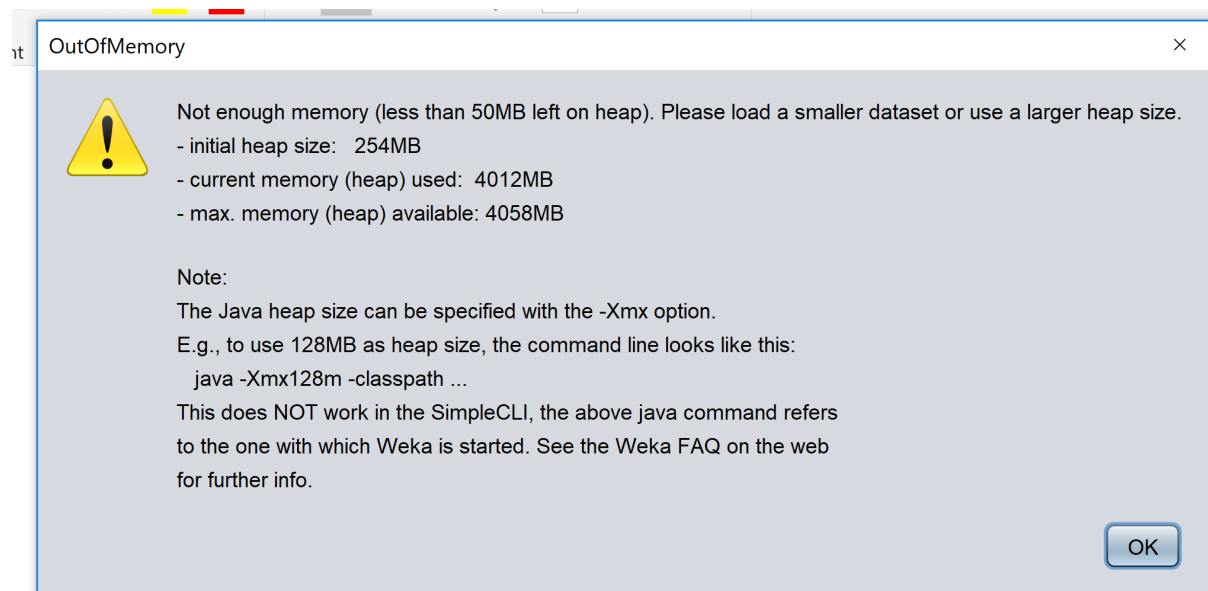


Figure 9. Running Apriori on supermarket2.



Figure 10. Running Apriori on supermarket1.