

Graduate Program in Software
CSIS 734-02: Data Mining & Predictive Analytics
 Assignment #2 (100 points)
Due Date: February 24th, 2018

Use the Apriori algorithm to mine association rules from the following transactional database. Let the support threshold be 60% (happened 3 times) and the confidence threshold be 100%.

1. How many times does the database need to be scanned?

The scan is used to find candidate itemsets. In this case, the database must be scanned twice.

2. What are the candidate itemsets for each database scan?

After scan one, the candidates are {Bread}, {Jelly}, {Milk}, {Peanut Butter} and {Beer}.

After scan two, the candidates are {Bread, Peanut Butter}.

3. What are the **final**** large (frequent) itemsets?**

The final large (frequent) itemset is taken from scan one results, after the support threshold is applied. It includes {Bread}, {Jelly}, {Milk}, {Peanut Butter} and {Beer}. Scan two on candidate 2-itemset {Bread, Peanut Butter} returns a support of only 2/5, or 40%. This is below the support threshold of 3/5, or 60%. Because it does not reach support threshold, no large 2-itemset is created.

4. What are the final association rules?

The final association rules are found for large itemsets. The last large itemset to be found was the first, containing {Bread} and {Peanut Butter}. The confidence for its variations are:

$$\begin{aligned} \text{Confidence (B} \rightarrow \text{P)} &= (B \vee P) / (B) \\ &= 2 / 4 \\ &= 50\% \end{aligned}$$

$$\begin{aligned} \text{Confidence (P} \rightarrow \text{B)} &= (B \vee P) / (P) \\ &= 2 / 2 \\ &= 100\% \end{aligned}$$

Key	TID	Itemset
B = Bread	1	B, J, P
J = Jelly	2	B, P
M = Milk	3	B, M, P
P = Peanut butter	4	R, B
R = Bread	5	R, M

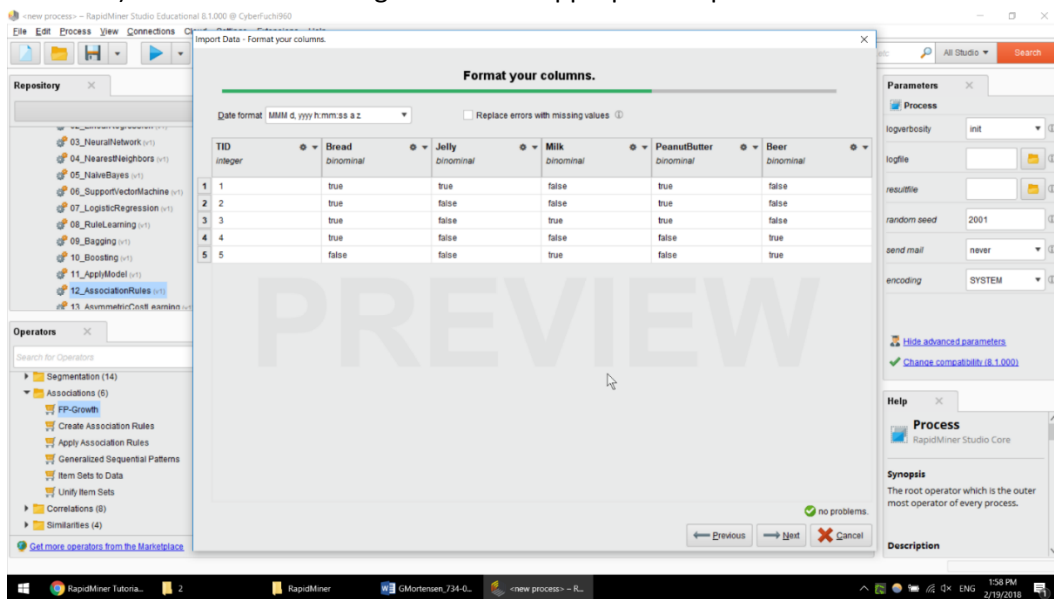
Scan ->	Candidate 1-itemset	Support	Large 1-itemset Support
	{B}	4	{B} 4
	{J}	1	{P} 3
	{M}	2	
	{P}	3	
	{R}	2	

Scan ->	Candidate 2-itemset	Support	Large 2-itemset Support
	{B, P}	2	---

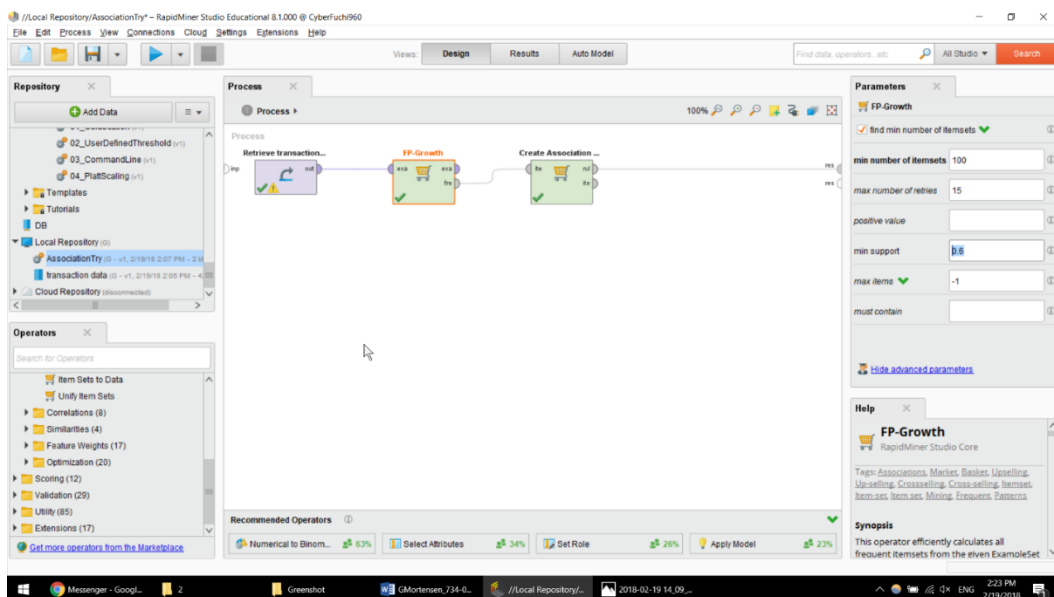
5. **(Extra Credits 10 points)** Do the following tasks: (1) select your favorite Data Mining tool, (2) convert the following table to the appropriate input format of the tool, (3) run the association rule mining process with the thresholds given in this assignment, (4) get a screen shot of the mining results, (5)** submit the screen shot with your answers to question 1 to 4, (6) please also indicate the name of your data mining tool.

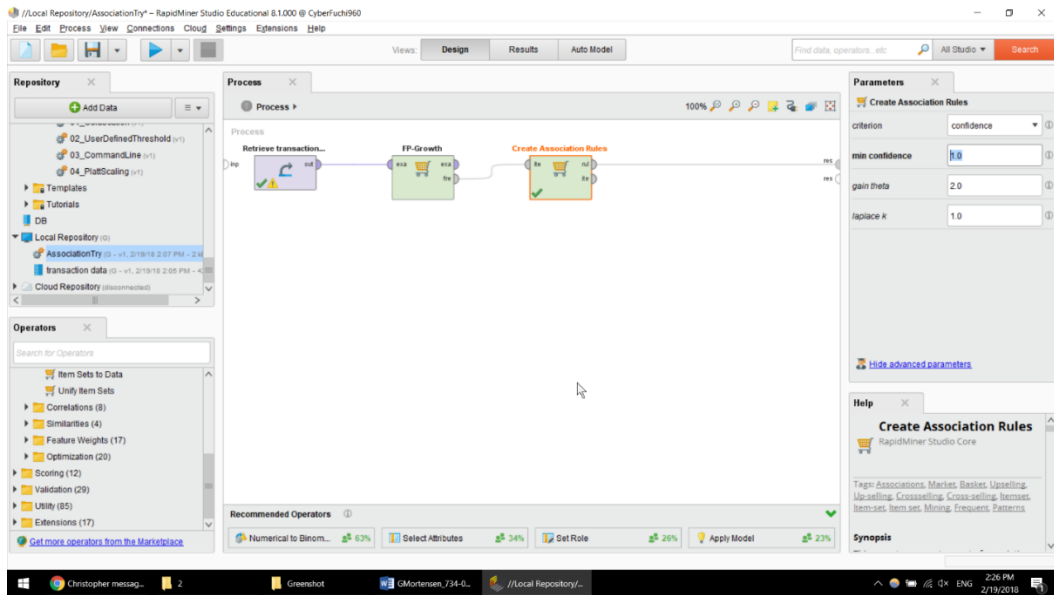
1) RapidMiner

2) convert the following table to the appropriate input format of the tool

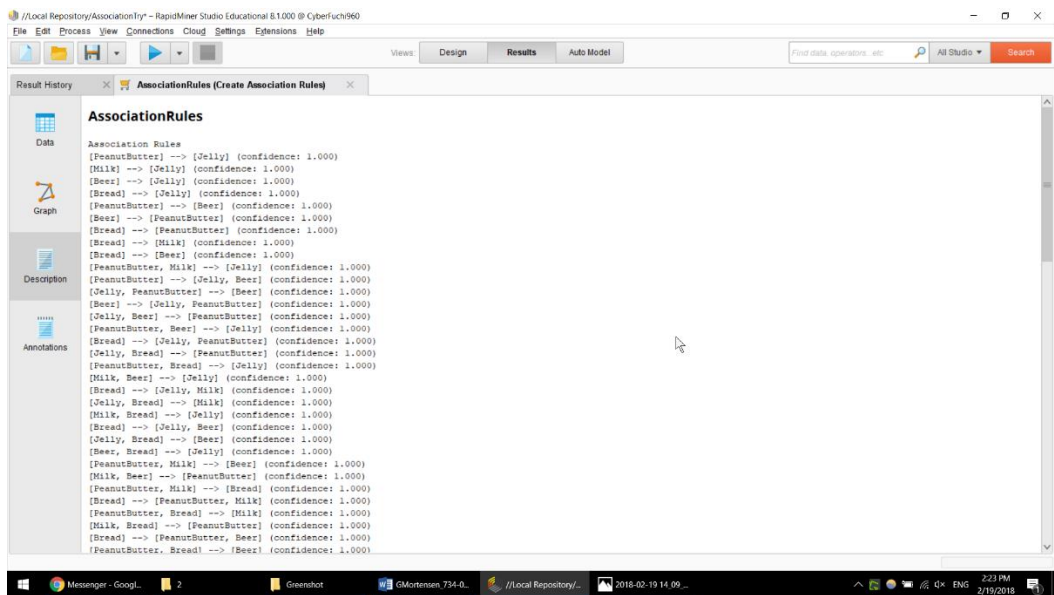


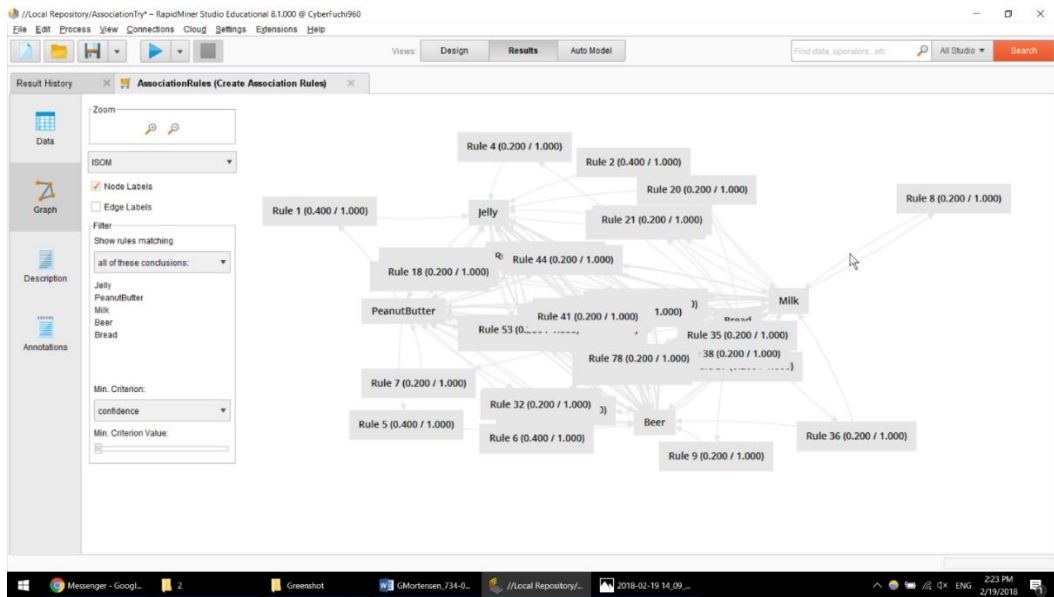
3) run the association rule mining process with the thresholds given in this assignment





4) get a screen shot of the mining results





Incomplete script with R:

This is an attempt at association analysis for data mining course.

it is also my first actual attempt at doing anything more than descriptive statistics with r.

INSTALL AND LOAD PACKAGES

pacman::p_load(arules, arulesViz)

DATA

Read transactional data from arules package

data("Groceries") # Load data

?Groceries # Help on data

str(Groceries) # Structure of data

summary(Groceries) # Includes 5 most frequent items

RULES

Set minimum support (minSup) to .001

Set minimum confidence (minConf) to .75

```
rules <- apriori(Groceries,  
  parameter = list(supp = 0.001, conf = 0.75))
```

options(digits=2)

inspect(rules[1:10])

PLOTS

Scatterplot of support x confidence (colored by lift)

plot(rules)

Graph of top 20 rules

```
plot(rules[1:20],  
  method = "graph",
```

```
control = list(type = "items"))
```

```
# Parallel coordinates plot of top 20 rules
```

```
plot(rules[1:20],
```

```
method = "paracoord",
```

```
control = list(reorder = TRUE))
```

```
# Matrix plot of antecedents and consequents
```

```
plot(rules[1:20],
```

```
method = "matrix",
```

```
control = list(reorder = TRUE))
```

```
# Grouped matrix plot of antecedents and consequents
```

```
plot(rules[1:20], method = "grouped")
```

Note: Assignment will be collected right before the class. **Don't forget to include your name on the top of your answer sheet.**