# 1   Introduction

**Mining association rules using ARS, from the SIPINA distribution.**

SIPINA is known for its decision tree induction algorithms. In fact, the distribution includes two other tools that are little known to the public: REGRESS, which is specialized in multiple linear regression, we described it in one of our tutorials[1]; and an association rules extraction tool, called simply Association Rule Software (**ARS**).

At that time (1997-1998), my idea was to build a set of independent executables organized around the same data management grid. Various tools for factor analysis and clustering, in addition to the tools which are included in the current distribution, were developed. But afterward, I realized that a functioning based on visual programming using streams of components was more advantageous in many respects, in terms of ease of use for the users, but also in terms of software evolution. TANAGRA was developed with this in mind.

This does not mean that the tools included in the SIPINA distribution are not interesting. In this tutorial, I describe the use of the ARS tool. Its interactivity with Excel spreadsheet is its main advantage. We launch the software from Excel using the "sipina.xla" add-in[2]. We can easily retrieve the rules in the spreadsheet. Then, we can explore them (the rules) using the Excel data handling capabilities. The ability to filter and sort rules according to different criteria is a great help in detecting interesting rules. This is a very important aspect because the profusion of rules can quickly confuse the data miner.

# 2   Dataset

The "**market_basket.xlsx**" data file describes the contents of **n = 1361** shopping carts (transactions). We have a library of **p = 303** products (items).

In average, each cart contains 9,5 products (min = 0, max = 303). And the products are purchased 42,7 times in average (min = 7, max = 167). The 5 most popular products are: "Eggs", "White bread", "2pct milk", "Potato chips" and "98pct fat free hamburger". The least popular ones are: "Celery", "Oats and Nuts Cereals", "Chicken legs", "Nasal spray" and "Daily Newspaper".

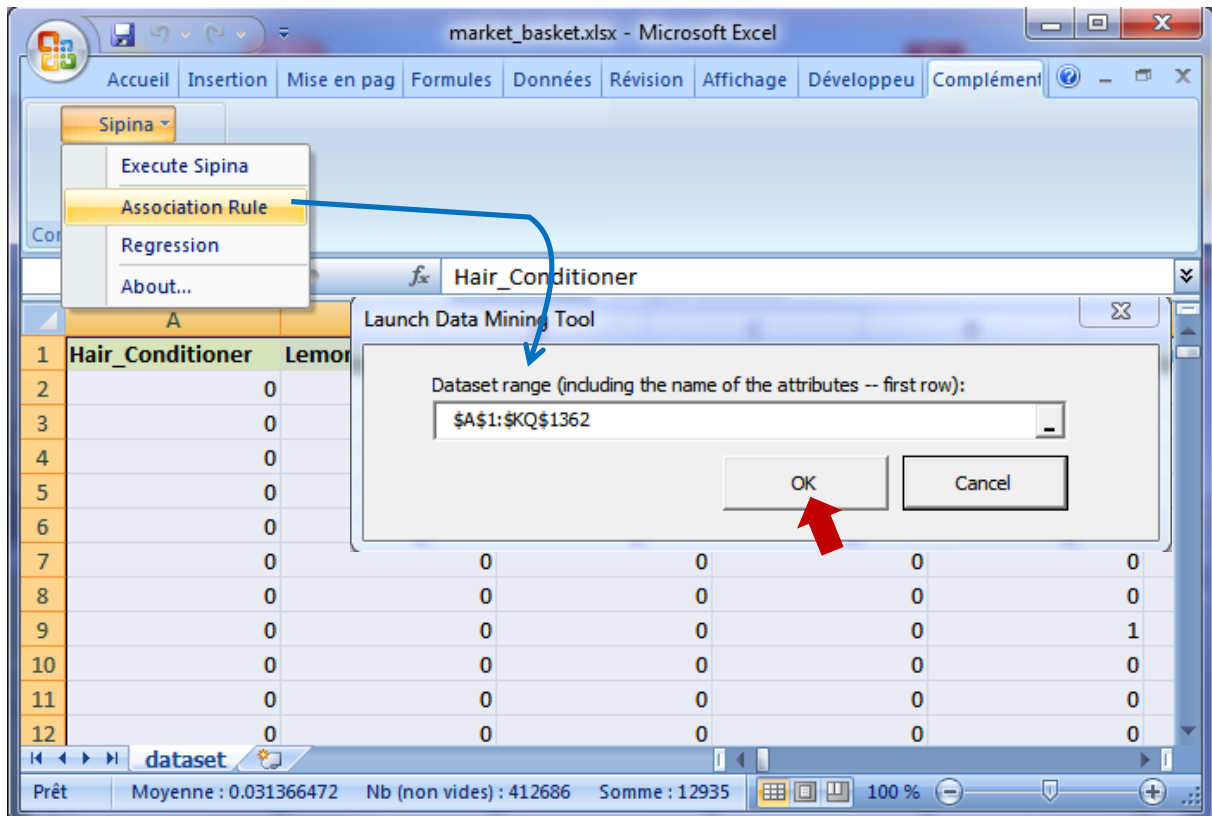# 3   Mining association rule with ARS

## 3.1   Data importation

We load the "**market_basket.xlsx**" data file into Excel (we must use Excel 2007 or later version because the number of items [p = 303] exceeds the limitation of 256 columns of Excel 2003 and earlier versions). **Note:** If you do not have the Excel 2007 version, the best solution is to import the "market_basket.txt" data file (tab separated text format) via the "FILE / OPEN / TEXT FILE FORMAT (*.txt)" menu[3]. This file is included in the archive that comes with this document.

---

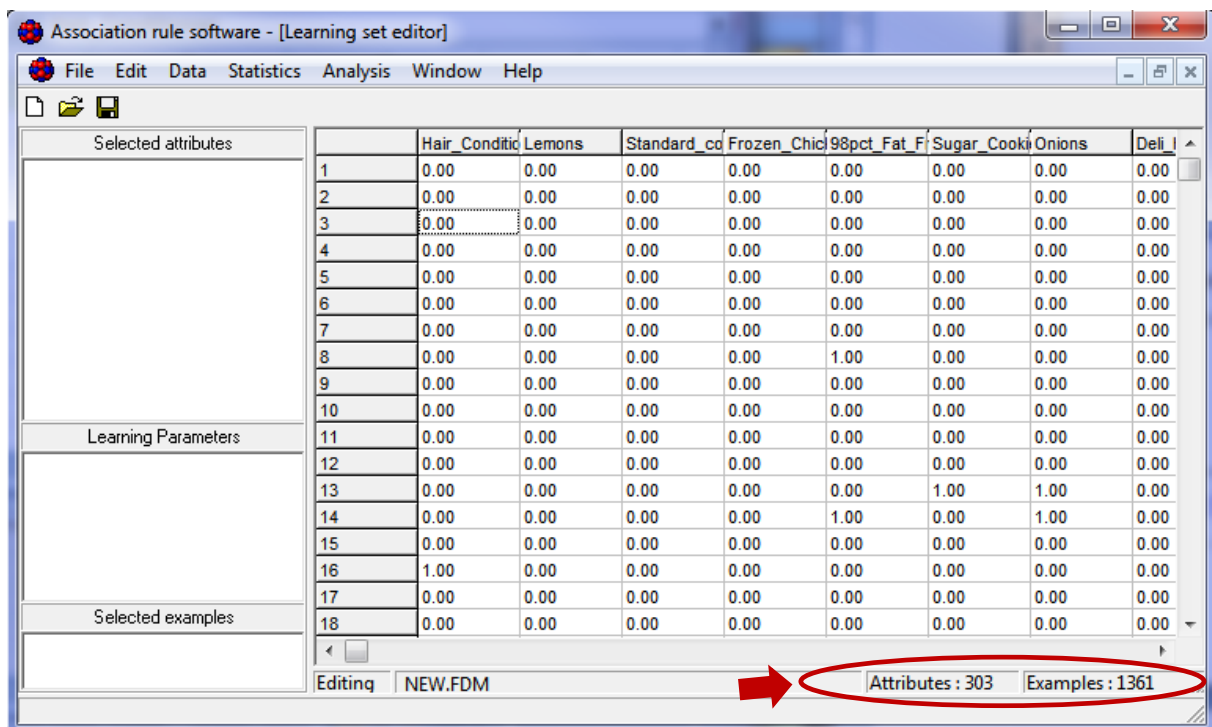[1] http://data-mining-tutorials.blogspot.fr/2011/08/regress-into-sipina-package.html

[2] http://data-mining-tutorials.blogspot.fr/2016/06/sipina-add-in-for-excel-2007-and-2010.html

[3] Sipina – Supported file format - http://data-mining-tutorials.blogspot.fr/2009/11/sipina-supported-file-format.html

---

After loading the data file into Excel, we select the data range and click the menu SIPINA / ASSOCIATION RULE. A setting dialog box appears. We check the coordinates of the data range and we click the OK button.
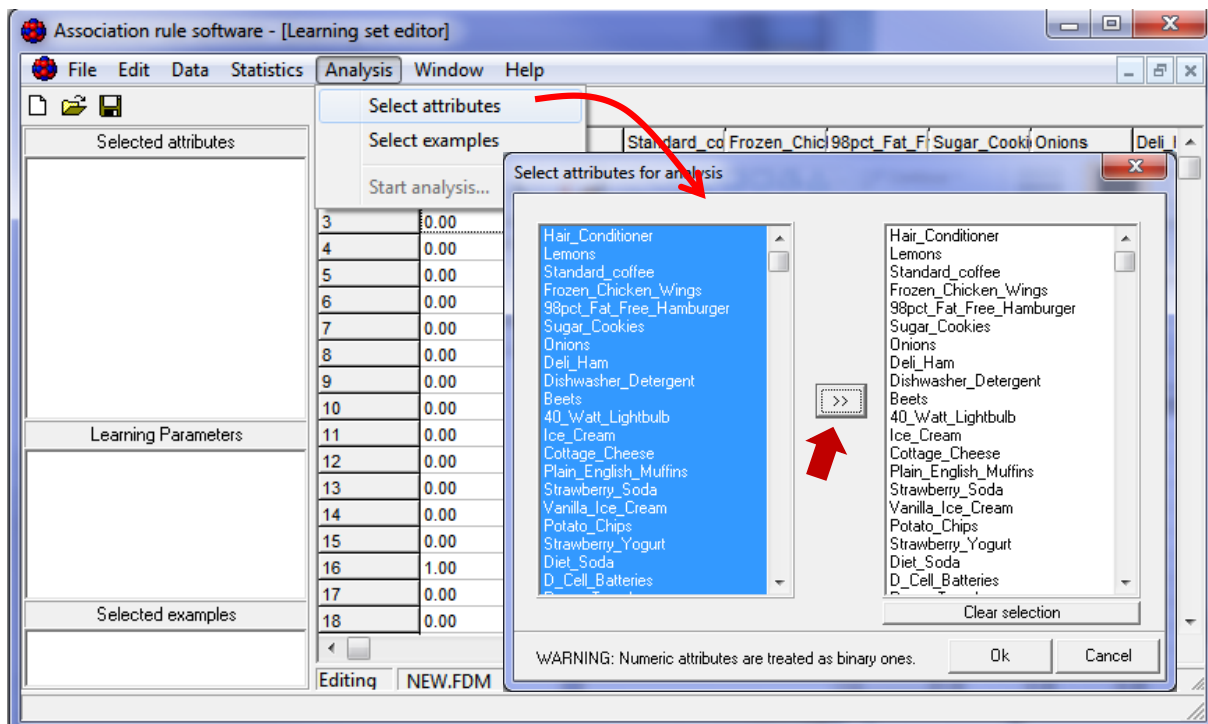


The software is automatically launched and the dataset is imported. We check that we have n = 1361 rows and p = 303 columns.
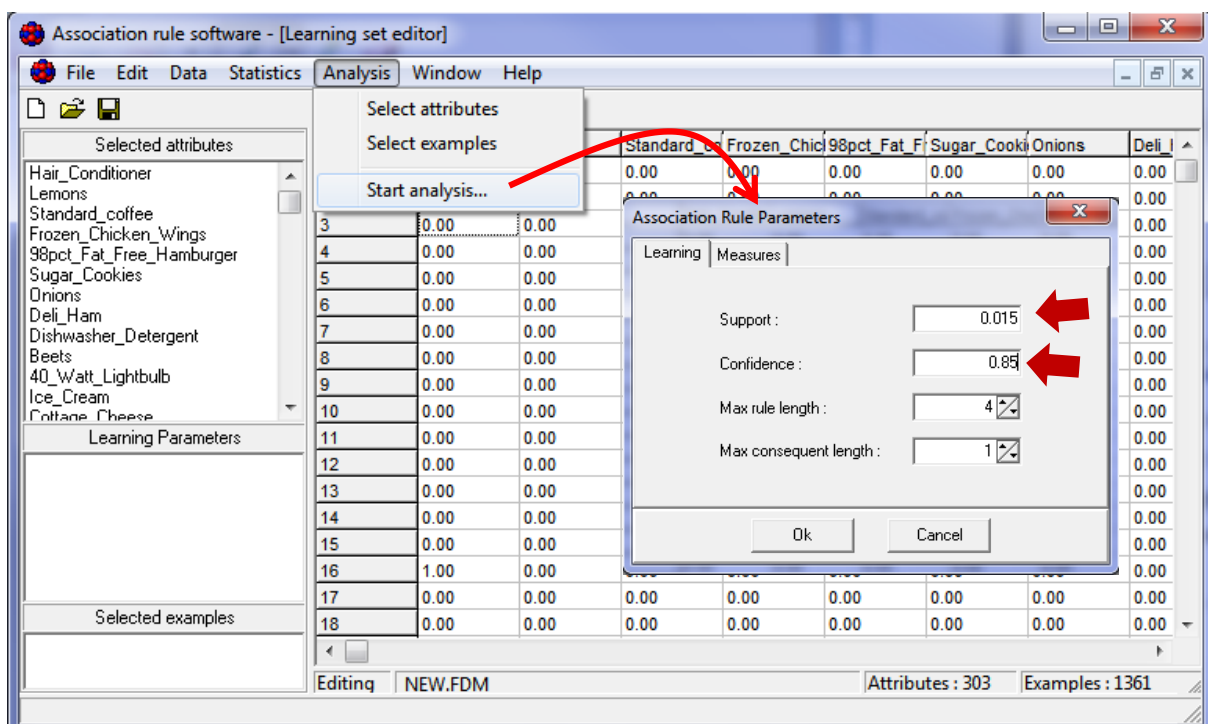
## 3.2    Selecting the items

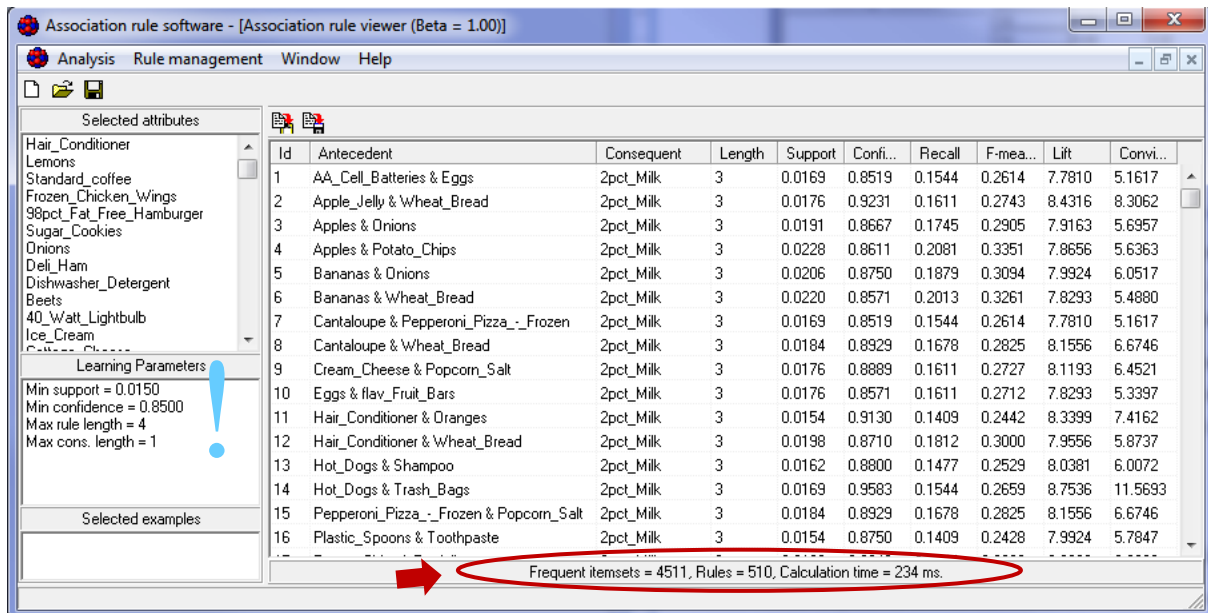To select the items to include into the analysis, we click the ANALYSIS / SELECT ATTRIBUTES menu. We select all the columns.



Because our variables are numerical, ARS takes them as dummy variables: 0, absence of the item in the shopping cart; 1 (it can be any value > 0), presence. When we process categorical variables, the tool automatically performs a dummy coding before starting processing.

## 3.3    Settings for the rule extraction process

To launch the association rule learning, we click the ANALYSIS / START ANALYSIS menu. A setting dialog box appears: we set the minimum support to 0.015 (we accept the rules which occurs at least 1361 x 0.015 ≈ 20 times into the database); the minimum confidence is 0.85; the maximum cardinal of the rule is 4 items; we have at most 1 item into the consequent of the rule.



When we validate these settings by clicking the OK button, we obtain a new window which contains the mined rules. Each rule is characterized by its antecedent, its consequent, and a set of numerical indicators (support, confidence, lift, etc.).[4]

We obtain **510 rules**. The first one (n°1) is:

**IF** purchase (AA_Cell_Batteries & Eggs) **THEN** purchase also (2pct_Milk)

The Support of the rule is:

P(AA_Cell_Batteries & Eggs & 2pct_Milk) = 0.0169

Its Confidence:

P(2pct_Milk / AA_Cell_Batteries & Eggs) = 0.8519

Its Lift:

$$\frac{P(2pct\_milk/AA\_Cell\_Batteries \ \& \ Eggs)}{P(2pct\_milk)} = 7.7810$$

Actually, nothing really distinguishes this tool from the association rule extractions components available in TANAGRA. My idea in this tutorial is to highlight in the next section the features... of Excel that allow us to examine the results (mined rules) better.
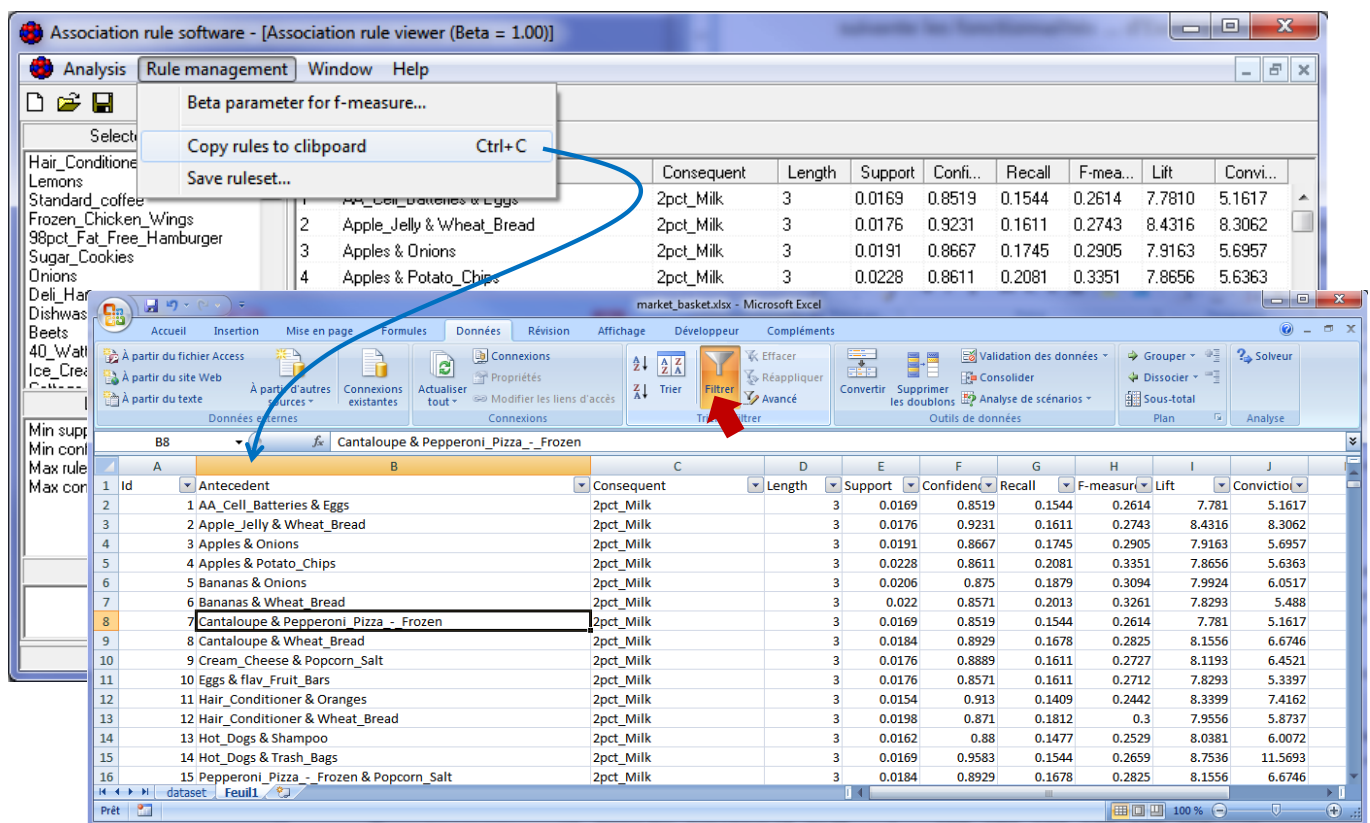
---

[4] "Interestingness measures for association rules" - http://data-mining-tutorials.blogspot.fr/2009/02/interestingness-measures-for.html

# 4   Exploring the rules under Excel

## 4.1   Retrieving the rules

Association rules learning algorithms often generate a large number of rules, 510 for our study. To be able to examine them, we must be able to organize the rules at our convenience: filter them according to several criteria, sort them according to the measures of interestingness, etc. The Excel's **Sort & Filter** tools are really appropriate for that purpose.

First, we must copy the rules into Excel. From ARS, we click the RULE MANAGEMENT / COPY RULES TO CLIPBOARD menu.



We create a new worksheet in Excel. We paste the rules. Then we activate the **FILTER** button (**Data** tab, **Sort & Filter** group)[5]. Arrows in columns headers enable to select various criteria for filtering the rules. The same functionalities are available in free spreadsheets such as Calc of LibreOffice[6] or Apache OpenOffice[7].

## 4.2   Sorting the rules according to a numerical indicator

We want to sort the rules according to the LIFT criterion. We click on the arrow of the corresponding column, we ask "Sort from the largest to the smallest".

---

[5] Our screenshots are based on the French version of Excel 2007.

[6] https://fr.libreoffice.org/

[7] https://www.openoffice.org
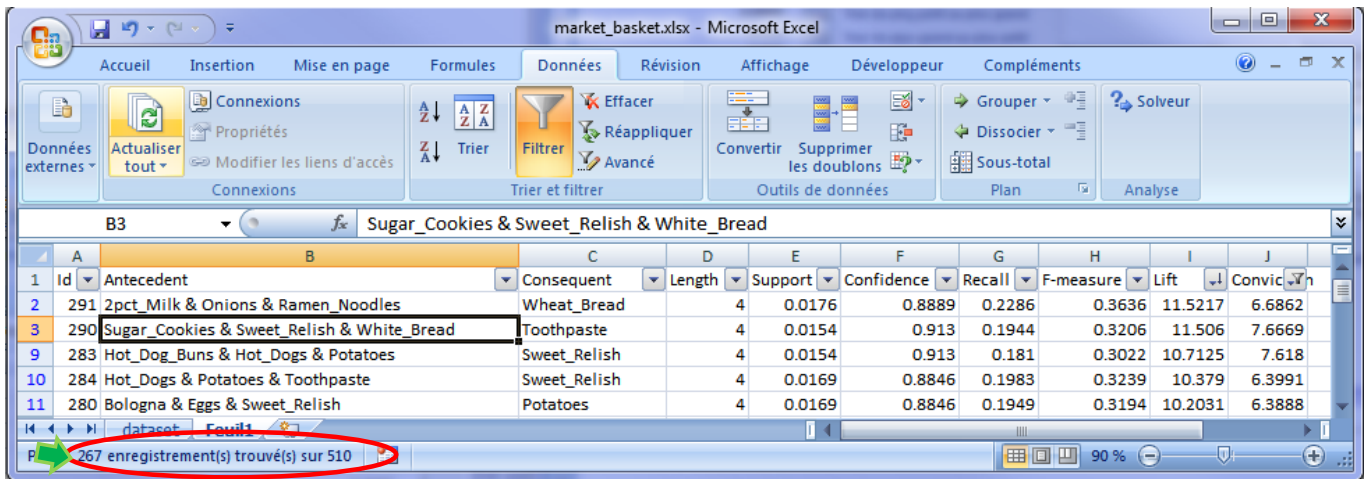
The rule with the highest LIFT is:

**IF** (2pct_Milk & Onions & Ramen_Noodles) **THEN** (Wheat_Bread) [**LIFT = 11.5217**]

## 4.3   Filtering according to a numerical condition

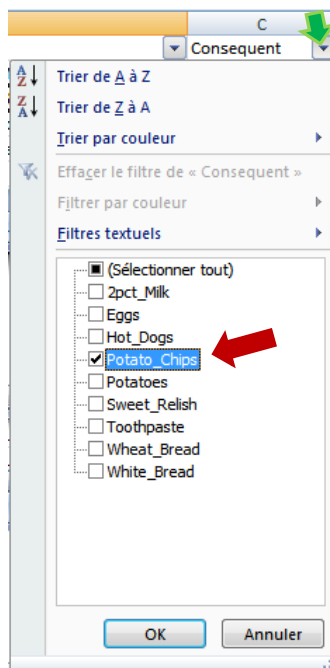We want to display only the rules with a conviction > 6.



To do this, we click on the corresponding arrow and we select the option "**Number Filters** / Greater than or Equal to…". We specify the threshold value 6.
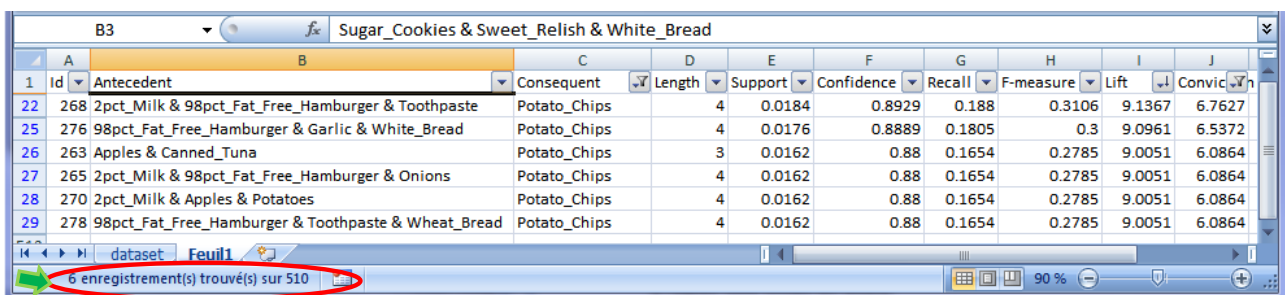
267 out of 510 rules meet this condition.

## 4.4 Filtering based on the consequent of the rule

**Among these rules,** we search the ones with the consequent "Purchase of (Potato_Chips)". We click the arrow in the header of the column and we select the corresponding item.
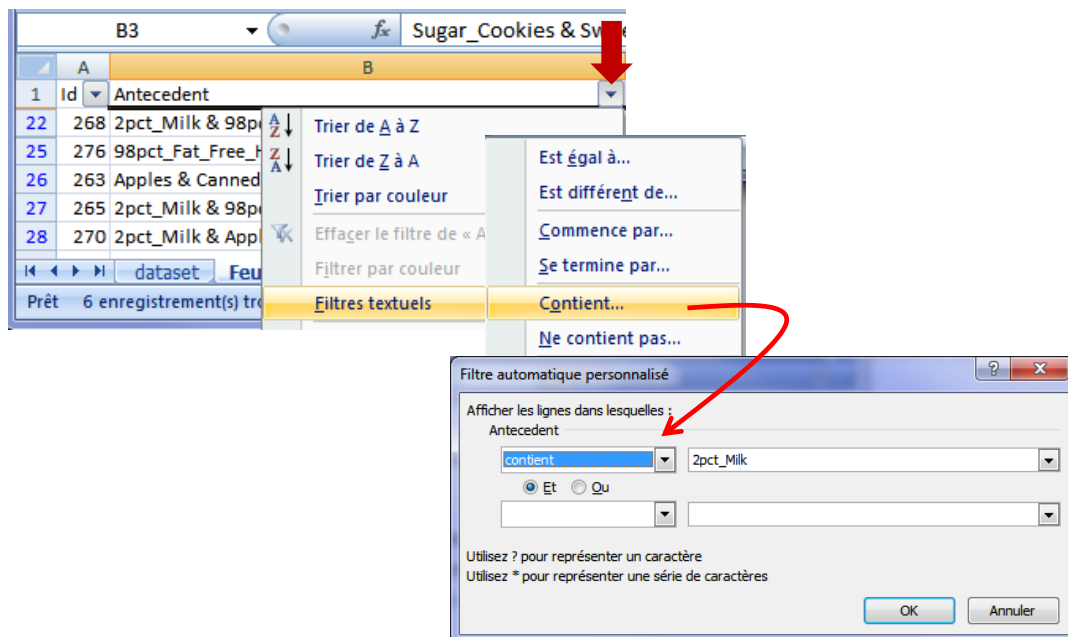


6 rules are highlighted.

## 4.1   Filtering based on the antecedent of the rule

**Among these rules**, we want to highlight the rules for which the antecedent contains the "2pct_Milk" item. We use the **Text Filters** tool, we set the condition.



3 rules are now displayed.



We can thus multiply the combinations to highlight the rules that best meet the specifications of our study. To be honest, this system is only really operational if we deal with a moderate-size rules base. But Excel is easy to use. It is one of the data miner favorite tools for a long time (KDnuggets Polls, Top Analytics / Data Science Tools, May 2017). We note in this tutorial that its functionalities are attractive for the processing of rule bases.

# 5   Conclusion

ARS is an academic software for association rule mining. It is included into the SIPINA distribution. In this tutorial, we describe its combined use with Excel. Among other things, the opportunities about the post-processing of the mined rules is particularly attractive.