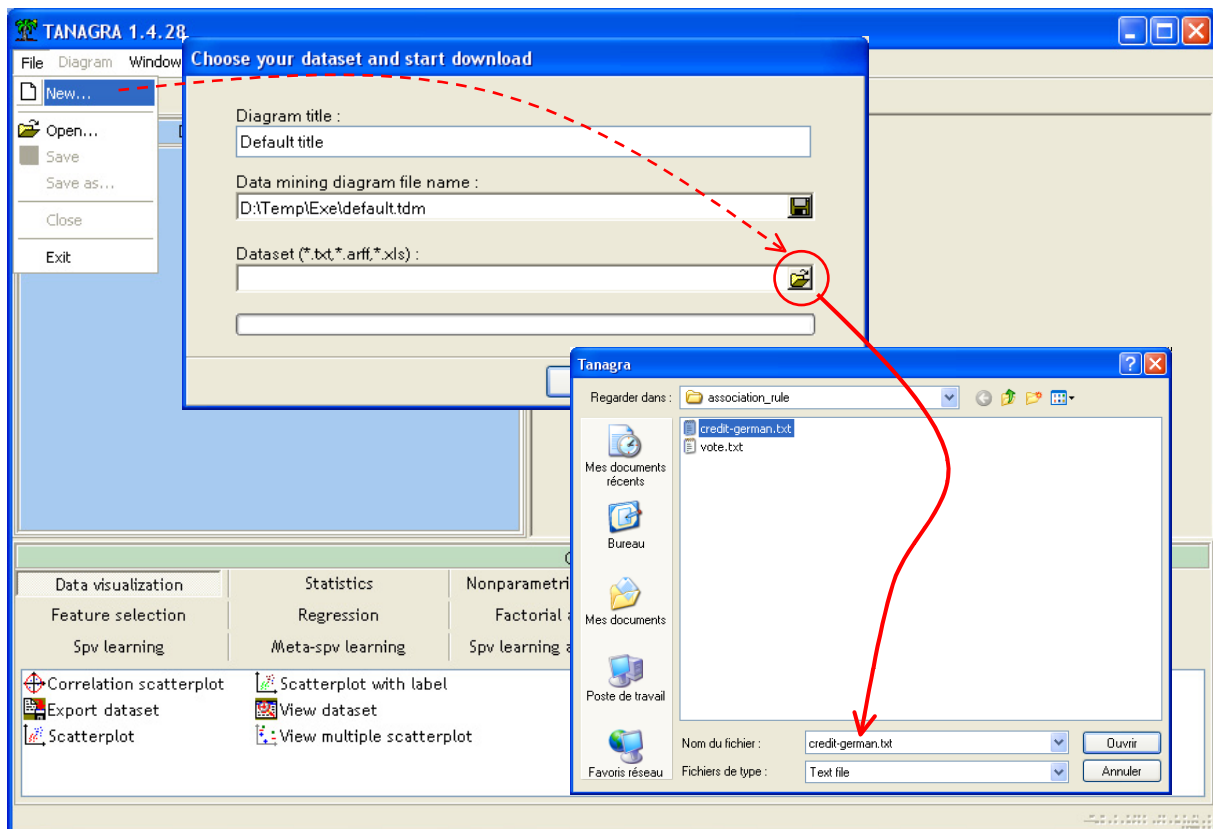


are: Tanagra 1.4.28, R 2.7.2 (arules package 0.6-6), Orange 1.0b2, RapidMiner Community Edition, Knime 1.3.5 and Weka 3.5.6. These programs load the data and perform the calculations in memory. When the size of the database increases, the real bottleneck is the memory available on our personal computer.

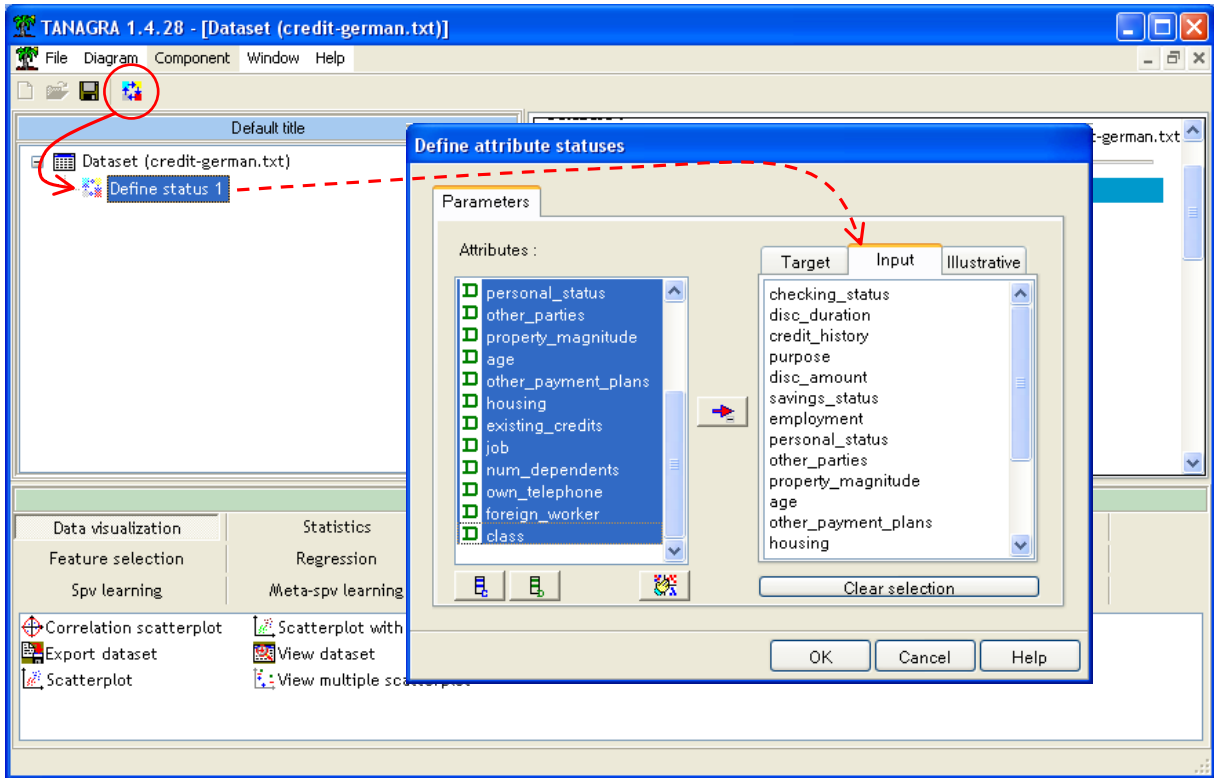
3 Tanagra (A Priori component)

Data importation and diagram initialization. After launching Tanagra, we activate the FILE / NEW menu. We select the CREDIT-GERMAN.TXT data file.

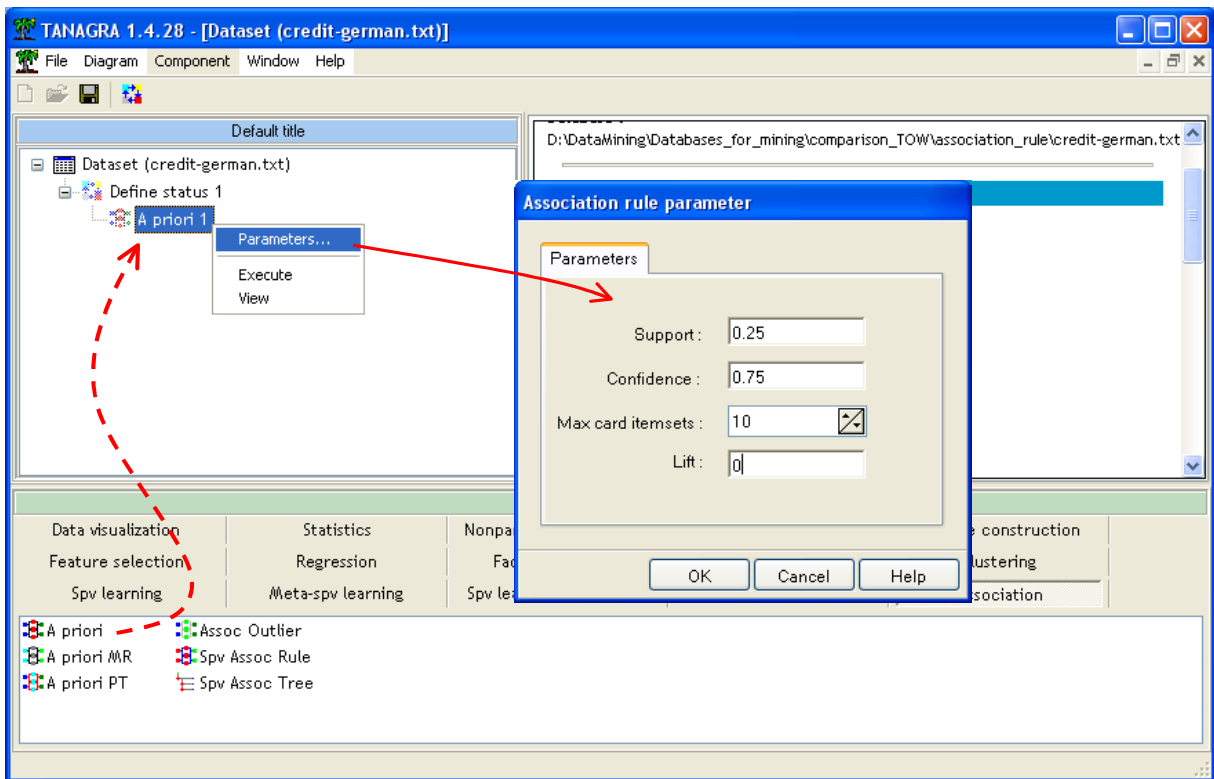


The dataset contains 19 variables and 1000 instances.

Specifying the status of the variables. In order to specify the status of each descriptor in the analysis, we use the DEFINE STATUS component. We set all the variables as INPUT.

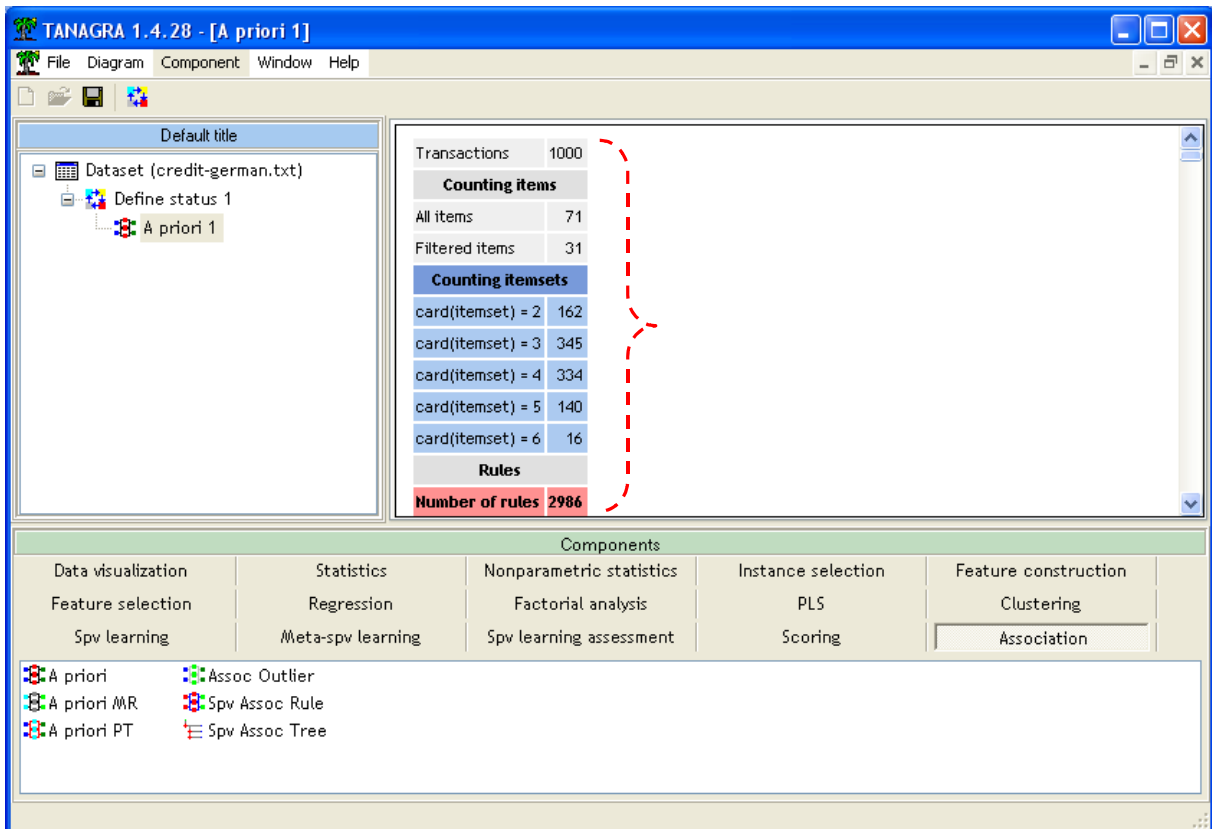


Extraction of association rules. We can insert the A PRIORI component now (ASSOCIATION tab). We activate the PARAMETERS menu in order to specify the parameters of the analysis.



The LIFT criterion is set to 0, thus it does not influence the extraction process. Our results will be directly comparable to the results of other tools which do not handle this parameter.

We click on VIEW menu.



Various indications are available: there are 71 items (attribute-value pair) into the dataset; 31 of them have a support ≥ 0.25 ; we see the number of itemsets of same length i.e. 162 itemsets with length = 2, etc.; thus, 2986 rules are extracted.

The rules are enumerated in the low part of the report. They are ranked according a decreasing value of the LIFT criterion.

Number of rules : 2986					
N°	Antecedent	Consequent	Lift	Support	Confidence
1	"other_payment_plans=none" - "existing_credits=one" - "own_telephone=none"	"credit_history=existing paid"	1.551	0.263	0.822
2	"num_dependents=one" - "class=good" - "credit_history=existing paid"	"other_payment_plans=none" - "existing_credits=one"	1.548	0.253	0.808
3	"other_parties=none" - "num_dependents=one" - "credit_history=existing paid"	"foreign_worker=yes" - "other_payment_plans=none" - "existing_credits=one"	1.534	0.319	0.769
4	"class=good" - "credit_history=existing paid"	"other_payment_plans=none" - "existing_credits=one"	1.534	0.289	0.801
5	"own_telephone=none" - "credit_history=existing paid"	"other_payment_plans=none" - "existing_credits=one"	1.527	0.263	0.797

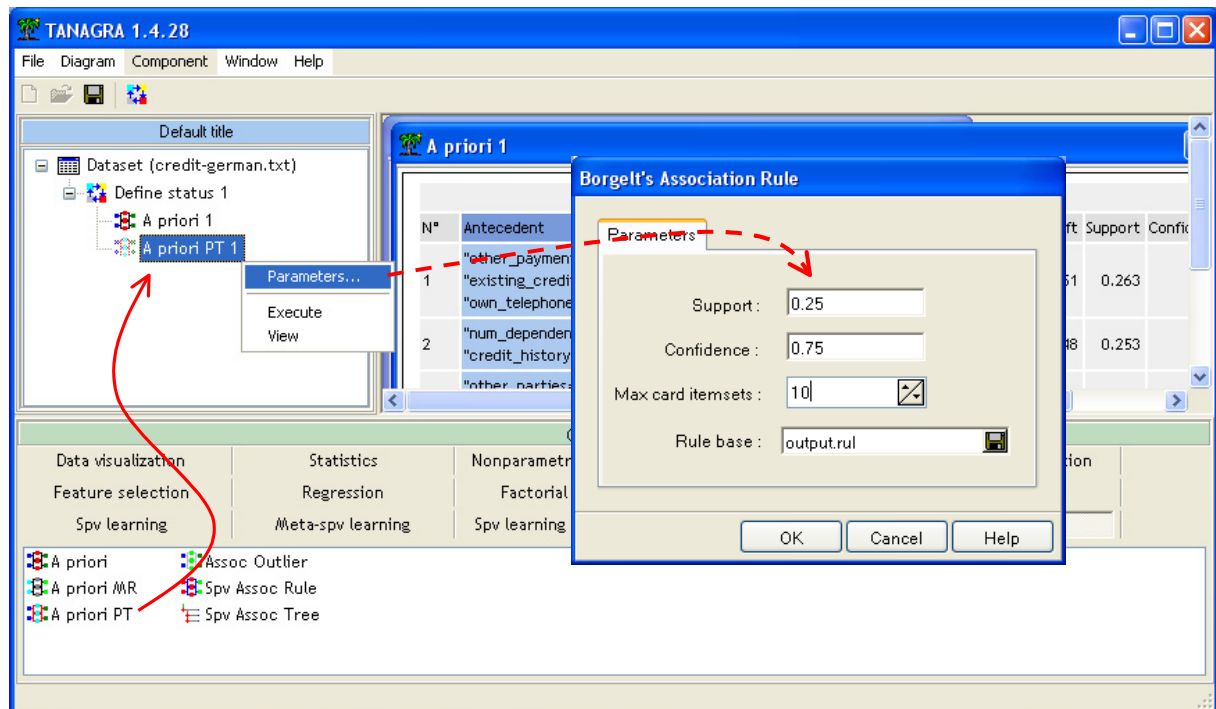


4 Tanagra (A Priori PT component)

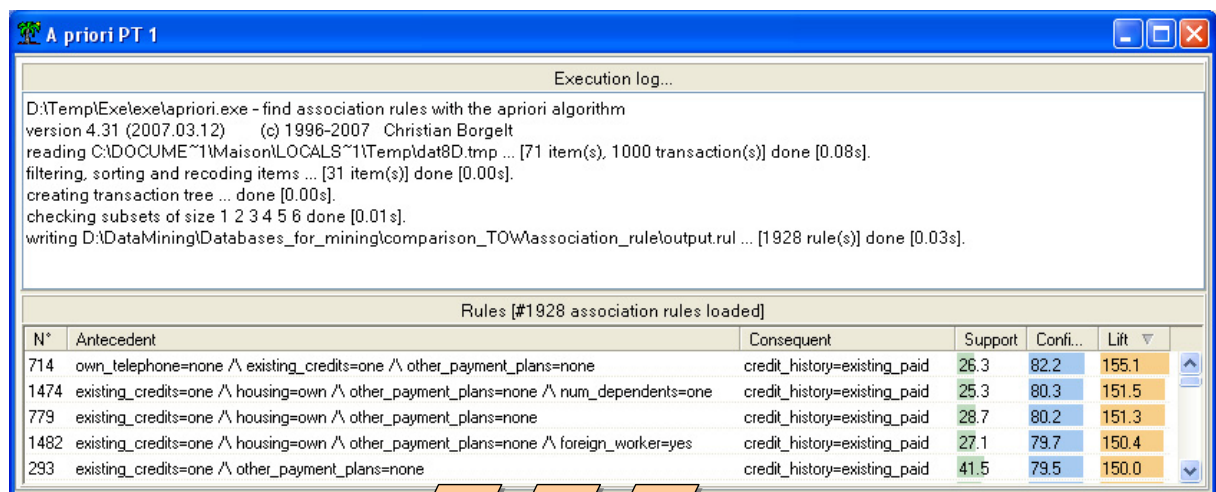
There is a second component dedicated to the association rule extraction in Tanagra. This is an external program "apriori.exe" of Christian BORGELT (<http://www.borgelt.net/apriori.html>). We have already written a first tutorial where we show the functioning of this tool previously (<http://data-mining-tutorials.blogspot.com/2008/11/association-rule-learning-using-prefix.html>).

BORGELT's program is really fast. But there is a slight limitation: only the rules with one item into the consequent can be generated. Therefore, for the same parameters above we obtain fewer rules.

In the previous diagram into Tanagra, we insert the A PRIORI PT component (ASSOCIATION tab). We set the parameters as follows.



We can click now on the VIEW menu.



In the upper part of the report, we can see the output of the BORGELT's program (version 4.31). In the lower part, the rules are enumerated. We can rank the rules according various numeric indicators by clicking on the column header. Because the component generates the rules with one item into the consequent, we obtain "only" 1928 rules.

5 R (arules package)

The « arules » package (<http://cran.univ-lyon1.fr/web/packages/arules/index.html>) allows extracting association rules with R (<http://www.r-project.org/>). This is also a version of the BORGELT's program, with the same limitation. In comparison with Tanagra, we must explicitly prepare the dataset before. The attribute-value representation must be transformed into a transactional data format. The operation is easy... if we read carefully the documentation.

Loading the « arules » package. We use the `library(.)` command in order to load the package

```
#charger le package
library(arules)
```

Data file importation and transformation. We import the dataset with the `read.table(.)` command, `summary(.)` gives some indications about the data characteristics.

```
#charger le fichier de données
setwd("D:/DataMining/Databases_for_mining/comparison_TOW/association_rule")
german <- read.table(file="credit-german.txt", header=T, dec=".", sep="\t")
summary(german)
```

We cannot extract rules from a data.frame, we must transform the internal format in "transactions".

```
#transformer les données attributs-variables
#en données transactionnelles
german.trans <- as(german, "transactions")
summary(german.trans)
```

We have always 71 items. R gives indications about the density of the dataset. We obtain a large number of rules if the density of the database is high.

```
> summary(german.trans)
transactions as itemMatrix in sparse format with
1000 rows (elements/itemsets/transactions) and
71 columns (items) and a density of 0.2676056

most frequent items:
      foreign_worker=yes      other_parties=none      num_dependents=one      other_payment_plans=none
              963              907              845              814
      housing=own              (Other)
              713              14758

element (itemset/transaction) length distribution:
sizes
19
1000

      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
      19      19      19      19      19      19

includes extended item information - examples:
      labels      variables      levels
1      checking_status=<0      checking_status      <0
2      checking_status=>=200      checking_status      >=200
3      checking_status=0<=X<200      checking_status      0<=X<200

includes extended transaction information - examples:
      transactionID
1              1
2              2
3              3
```

Extraction of rules. The following instructions extract the association rules.

```
#extraction des règles
german.regles <- apriori(german.trans,parameter=
      list(supp=0.25,conf=0.75,minlen=2,maxlen=10,target="rules"))
summary(german.regles)
```

We obtain again the BORGELT's program output.

```
> #extraction des règles
> german.regles <- apriori(german.trans,parameter=
+       list(supp=0.25,conf=0.75,minlen=2,maxlen=10,target="rules"))

parameter specification:
confidence minval smax arem  aval originalSupport support minlen maxlen target  ext
      0.75   0.1   1 none FALSE      TRUE   0.25     2    10 rules FALSE

algorithmic control:
filter tree heap memopt load sort verbose
      0.1 TRUE TRUE  FALSE TRUE     2    TRUE

apriori - find association rules with the apriori algorithm
version 4.21 (2004.05.09)      (c) 1996-2004  Christian Borgelt
set item appearances ...[0 item(s)] done [0.00s].
set transactions ...[71 item(s), 1000 transaction(s)] done [0.00s].
sorting and recoding items ... [31 item(s)] done [0.00s].
creating transaction tree ... done [0.00s].
checking subsets of size 1 2 3 4 5 6 done [0.01s].
writing ... [1928 rule(s)] done [0.00s].
creating S4 object ... done [0.00s].
```

Some characteristics of the generated rule base are also available.

```
> summary(german.regles)
set of 1928 rules

rule length distribution (lhs + rhs):sizes
  2  3  4  5  6
112 487 783 476 70

      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
      2.000  3.000  4.000  3.951  5.000  6.000

summary of quality measures:
      support      confidence      lift
Min.   :0.2500   Min.   :0.7500   Min.   :0.8894
1st Qu.:0.2700   1st Qu.:0.8397   1st Qu.:0.9972
Median :0.2990   Median :0.8906   Median :1.0121
Mean   :0.3294   Mean   :0.8859   Mean   :1.0437
3rd Qu.:0.3580   3rd Qu.:0.9485   3rd Qu.:1.0475
Max.   :0.8800   Max.   :0.9937   Max.   :1.5507

mining info:
      data ntransactions support confidence
german.trans      1000      0.25      0.75
```

Visualization of the rules. The inspect(.) command enables to visualize the details of rules. We show only the first 10 rules here.

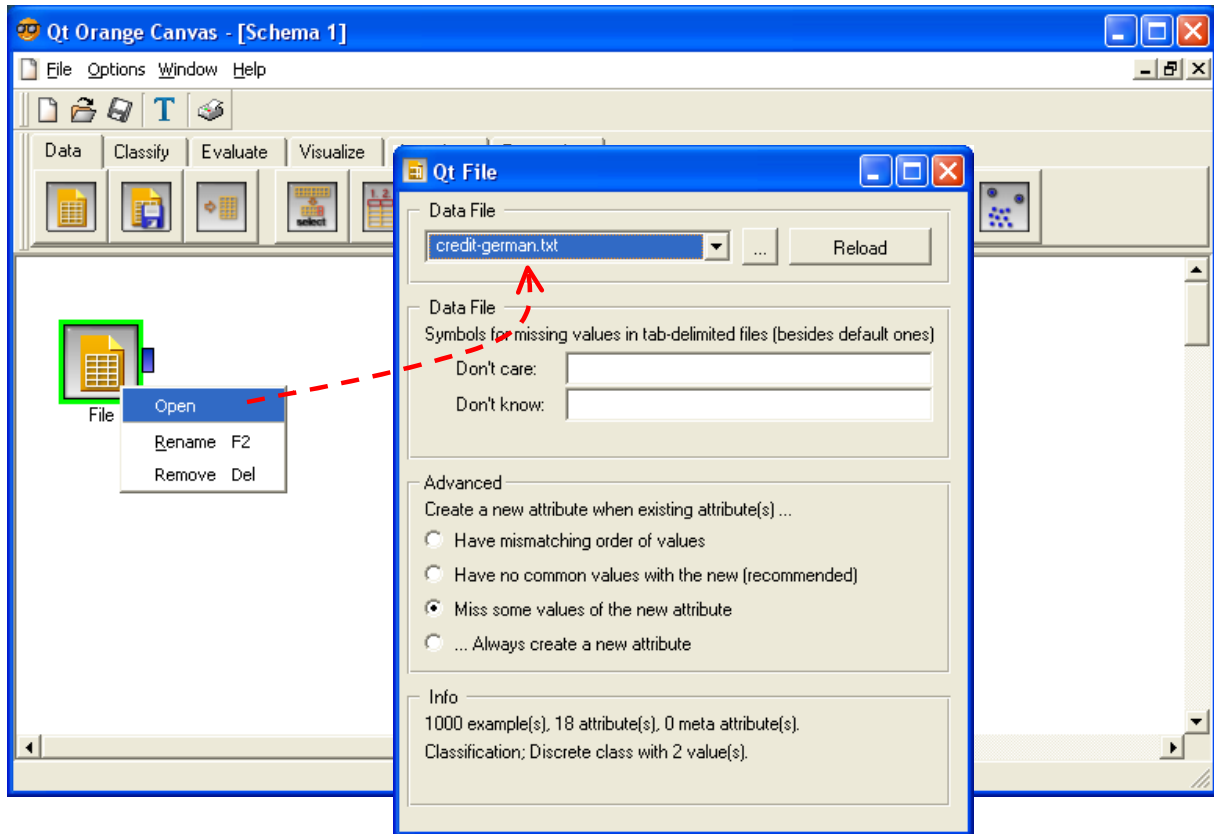
```
> #afficher les 10 premières règles trouvées
> inspect(german.regles[1:10])
lhs                                rhs                                support confidence lift
1 {checking_status=0<=X<200}      => {foreign_worker=yes}           0.264  0.9814126 1.0191201
2 {checking_status=<0}             => {foreign_worker=yes}           0.259  0.9452555 0.9815737
3 {purpose=radio/tv}              => {num_dependents=one}           0.250  0.8928571 1.0566357
4 {purpose=radio/tv}              => {foreign_worker=yes}           0.275  0.9821429 1.0198784
5 {property_magnitude=real estate} => {foreign_worker=yes}           0.262  0.9290780 0.9647747
6 {credit_history=critical/other existing } => {other_payment_plans=none}    0.251  0.8566553 1.0524021
7 {credit_history=critical/other existing } => {other_parties=none}           0.268  0.9146758 1.0084628
8 {credit_history=critical/other existing } => {foreign_worker=yes}           0.279  0.9522184 0.9888042
9 {class=bad}                      => {num_dependents=one}           0.254  0.8466667 1.0019724
10 {class=bad}                    => {other_parties=none}           0.272  0.9066667 0.9996325
```

We can also rank the rules according to a rule quality indicator. We show here the first 5 rules according to the LIFT criterion. We obtain the same rules as the A PRIORI PT component of Tanagra.

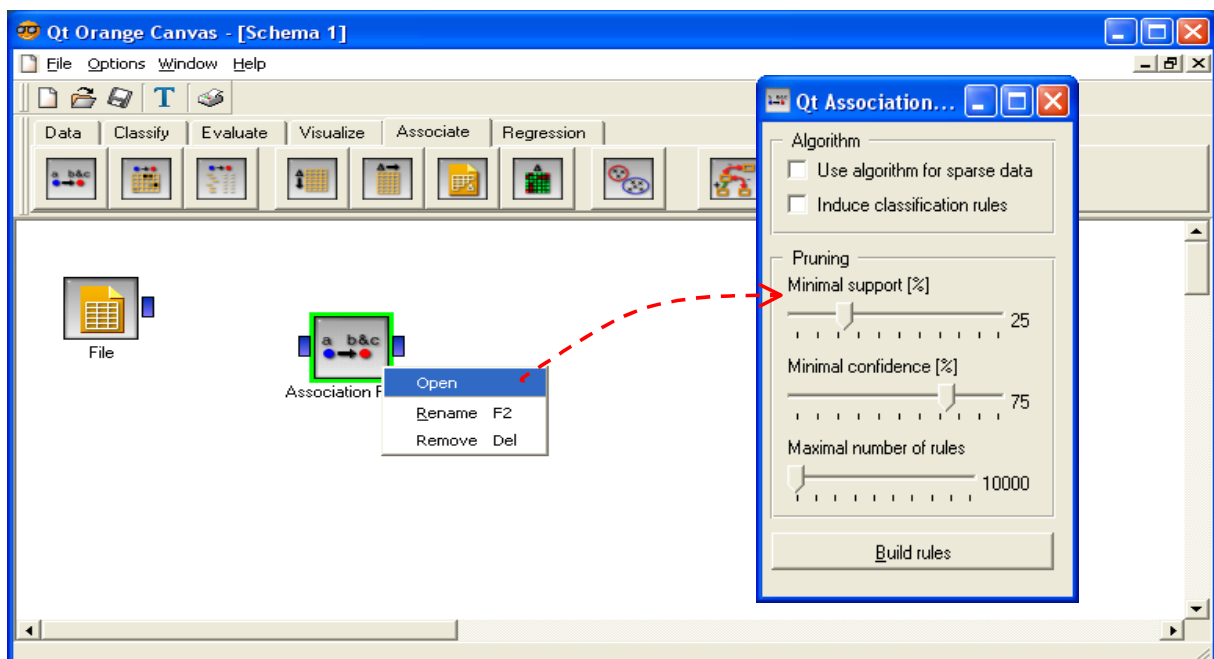
```
> #afficher les 5 règles avec le lift le + élevé
> regles.triees <- sort(german.regles,by="lift")
> inspect(regles.triees[1:5])
lhs                                rhs                                support confidence lift
1 {other_payment_plans=none,
  existing_credits=one,
  own_telephone=none}              => {credit_history=existing paid}    0.263  0.8218750 1.550708
2 {other_payment_plans=none,
  housing=own,
  existing_credits=one,
  num_dependents=one}              => {credit_history=existing paid}    0.253  0.8031746 1.515424
3 {other_payment_plans=none,
  housing=own,
  existing_credits=one}            => {credit_history=existing paid}    0.287  0.8016760 1.512596
4 {other_payment_plans=none,
  housing=own,
  existing_credits=one,
  foreign_worker=yes}              => {credit_history=existing paid}    0.271  0.7970588 1.503885
5 {other_payment_plans=none,
  existing_credits=one}            => {credit_history=existing paid}    0.415  0.7950192 1.500036
```


6 Orange

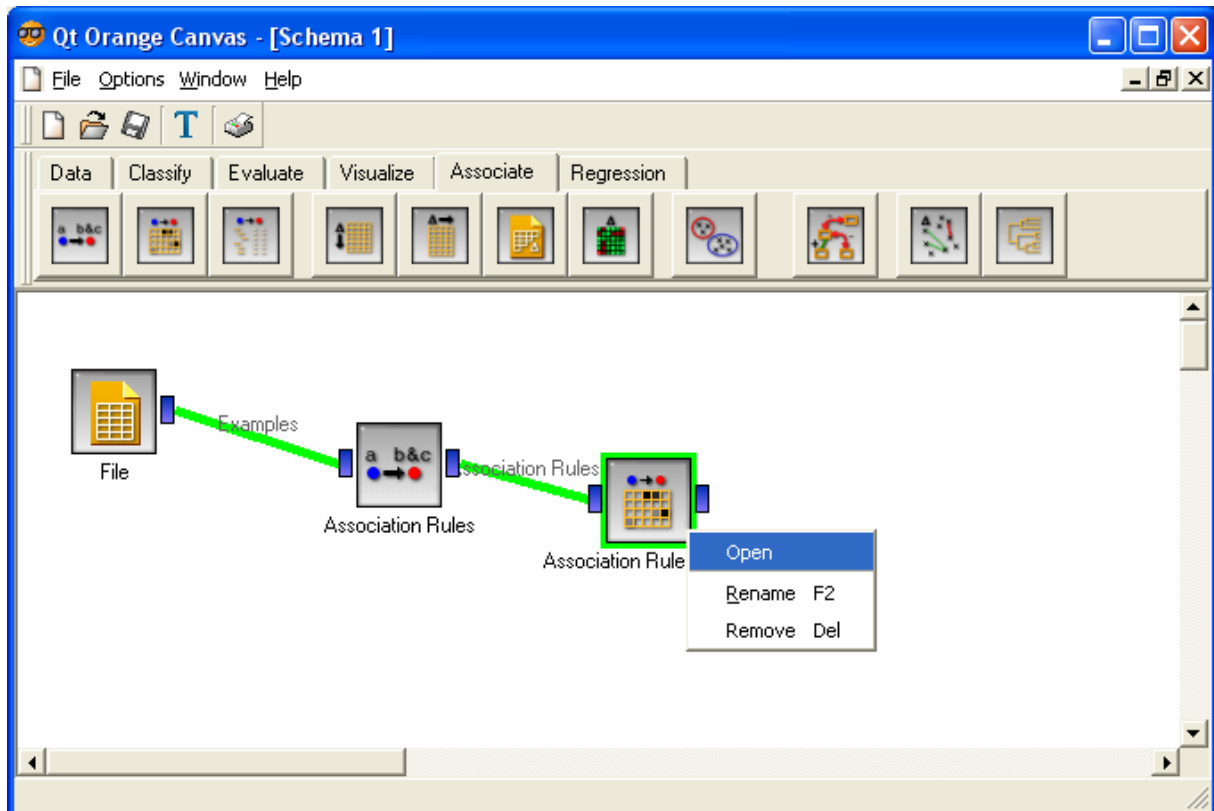
Creation of a "schema" and data importation. When we launch Orange, a new empty schema is available. We add the FILE component (DATA tab). We select the data file.



Rule extraction and visualization. We insert the ASSOCIATION RULES component (ASSOCIATE tab). We set the parameters of our analysis.



Then we add the ASSOCIATION RULES VIEWER component, we connect the components. We can now set the connection between FILE and ASSOCIATION RULES. The calculation is launched.

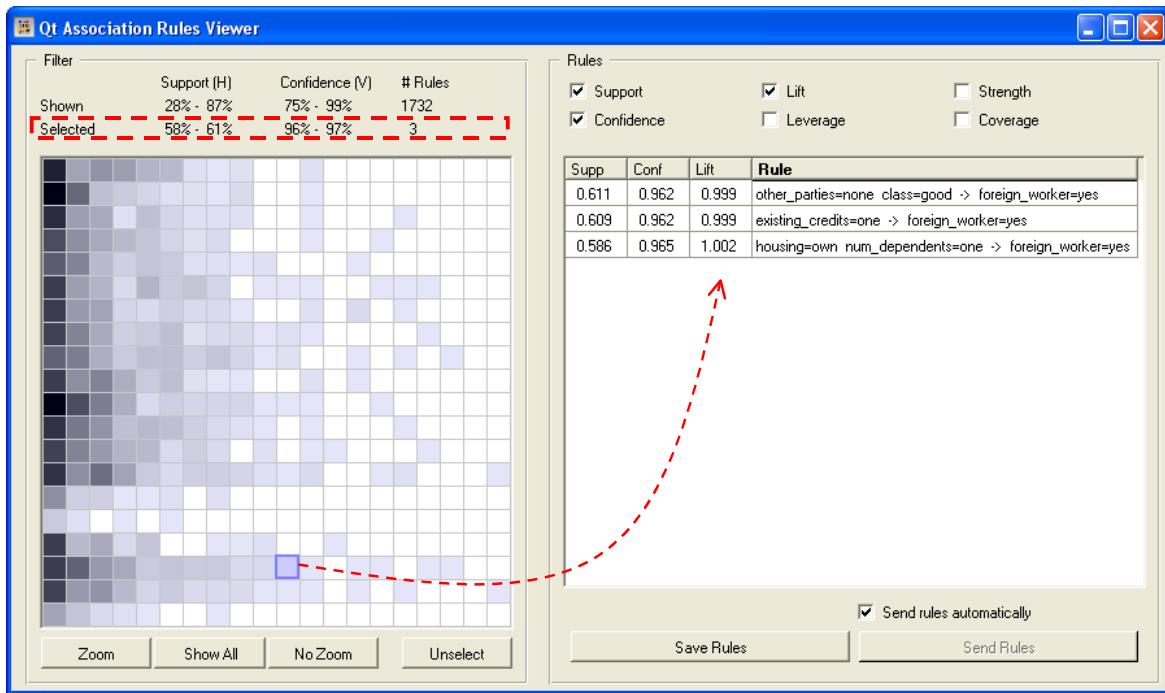


To see the rules, we click on the VIEW menu of the ASSOCIATION RULES VIEWER component.

The screenshot shows the Qt Association Rules Viewer window. It features a 'Filter' section with a heatmap visualization and a 'Rules' section with a table of association rules. The 'Rules' section includes checkboxes for 'Support', 'Confidence', 'Lift', 'Strength', 'Leverage', and 'Coverage'. The 'Send rules automatically' checkbox is checked. The 'Save Rules' and 'Send Rules' buttons are visible at the bottom.

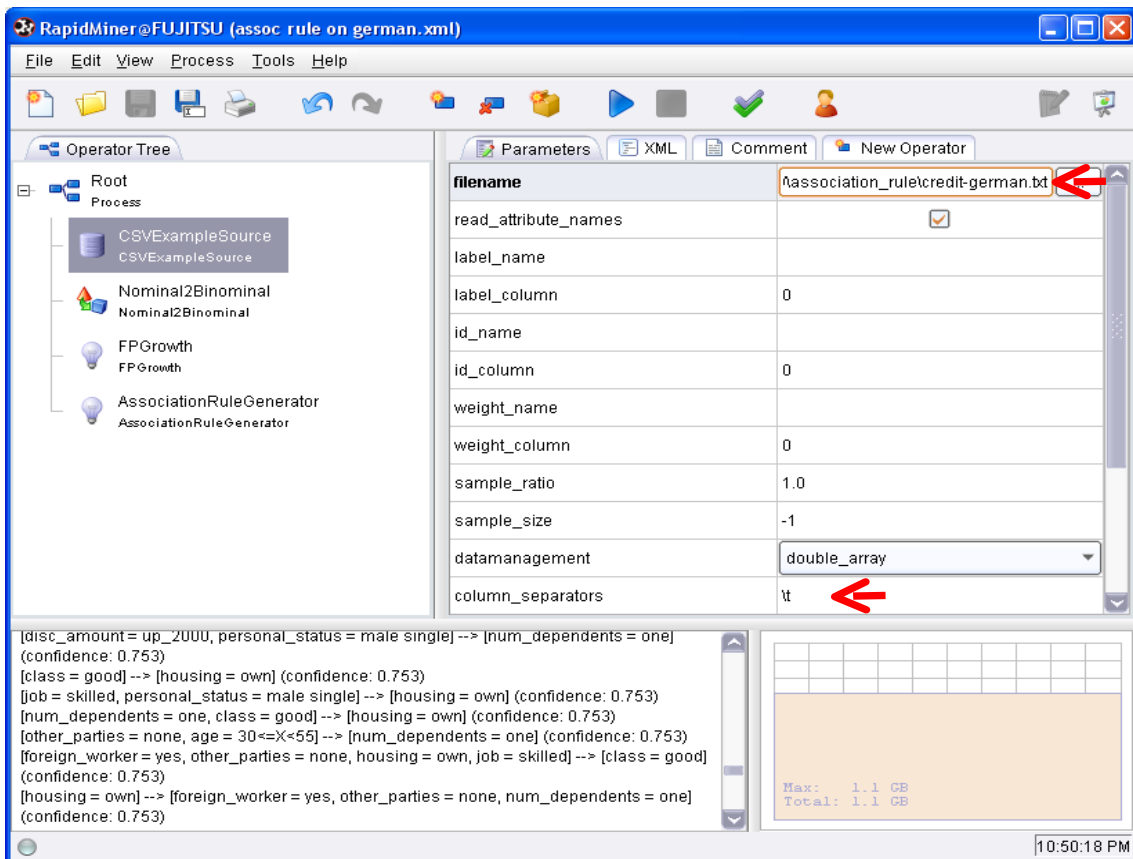
Supp	Conf	Lift	Rule
0.289	0.801	1.534	credit_history=existing paid class=good -> other_paymer
0.319	0.769	1.534	credit_history=existing paid other_parties=none num_dep
0.330	0.795	1.523	credit_history=existing paid other_parties=none num_dep
0.319	0.790	1.513	credit_history=existing paid other_parties=none num_dep
0.365	0.790	1.513	credit_history=existing paid num_dependents=one -> otl
0.350	0.758	1.512	credit_history=existing paid num_dependents=one -> otl
0.356	0.754	1.505	credit_history=existing paid other_parties=none -> other
0.350	0.785	1.503	credit_history=existing paid num_dependents=one foreign
0.370	0.784	1.502	credit_history=existing paid other_parties=none -> other
0.415	0.795	1.500	other_payment_plans=none existing_credits=one -> crec
0.415	0.783	1.500	credit_history=existing paid -> other_payment_plans=nor
0.365	0.793	1.497	other_payment_plans=none existing_credits=one num_de
0.350	0.761	1.492	other_payment_plans=none existing_credits=one num_de
0.396	0.790	1.491	other_payment_plans=none existing_credits=one foreign
0.319	0.820	1.491	credit_history=existing paid other_parties=none other_pa
0.319	0.792	1.491	credit_history=existing paid other_parties=none other pa

The visualization window is really original. We can graphically select a group of rules according to a range of the value of two criteria. Various criteria can be used; we can also rank the rules here.

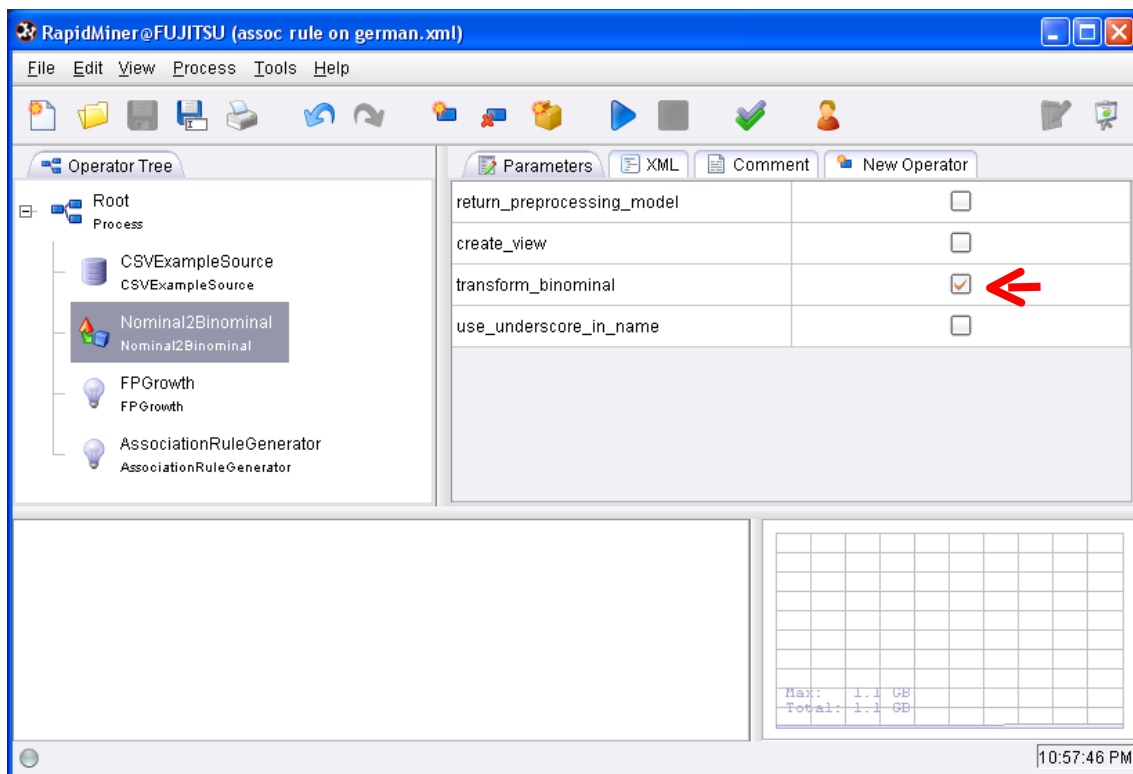


7 RapidMiner

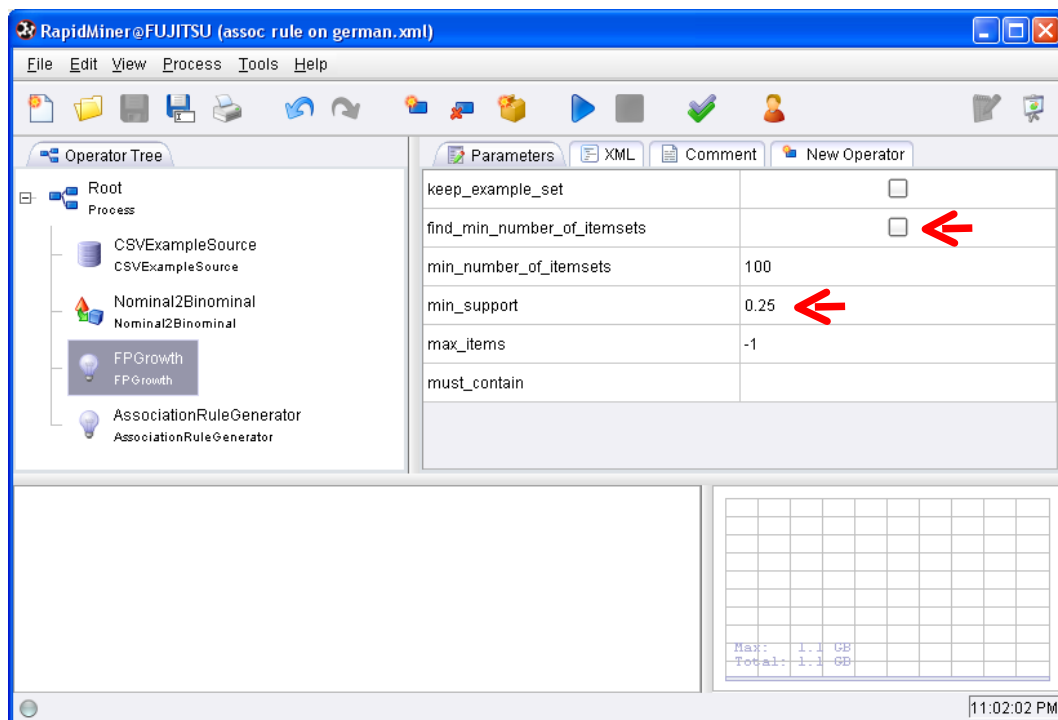
Creating a new operator tree. With RapidMiner, it is more convenient to define before all the sequences of operations before starting the computation. The first operator is **CSVEXAMPLESOURCE**, we set the file name and the column separator.



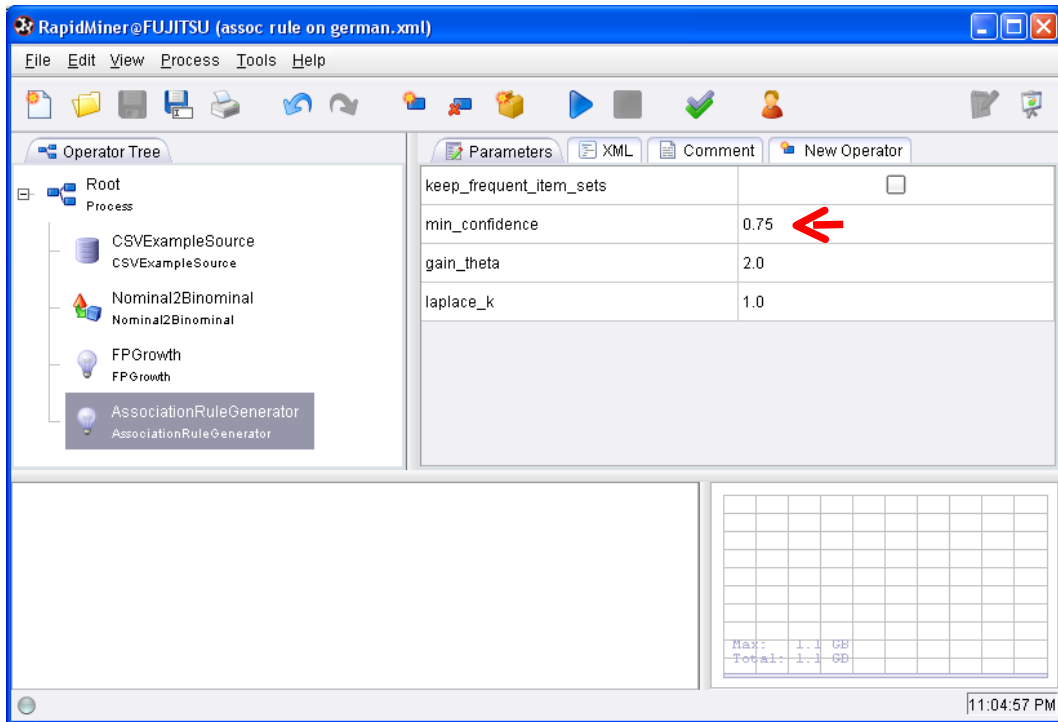
RapidMiner cannot create association rules from an attribute-value dataset. We must recode the variable into a set of binary columns with the **NOMINAL2BINOMIAL** component.



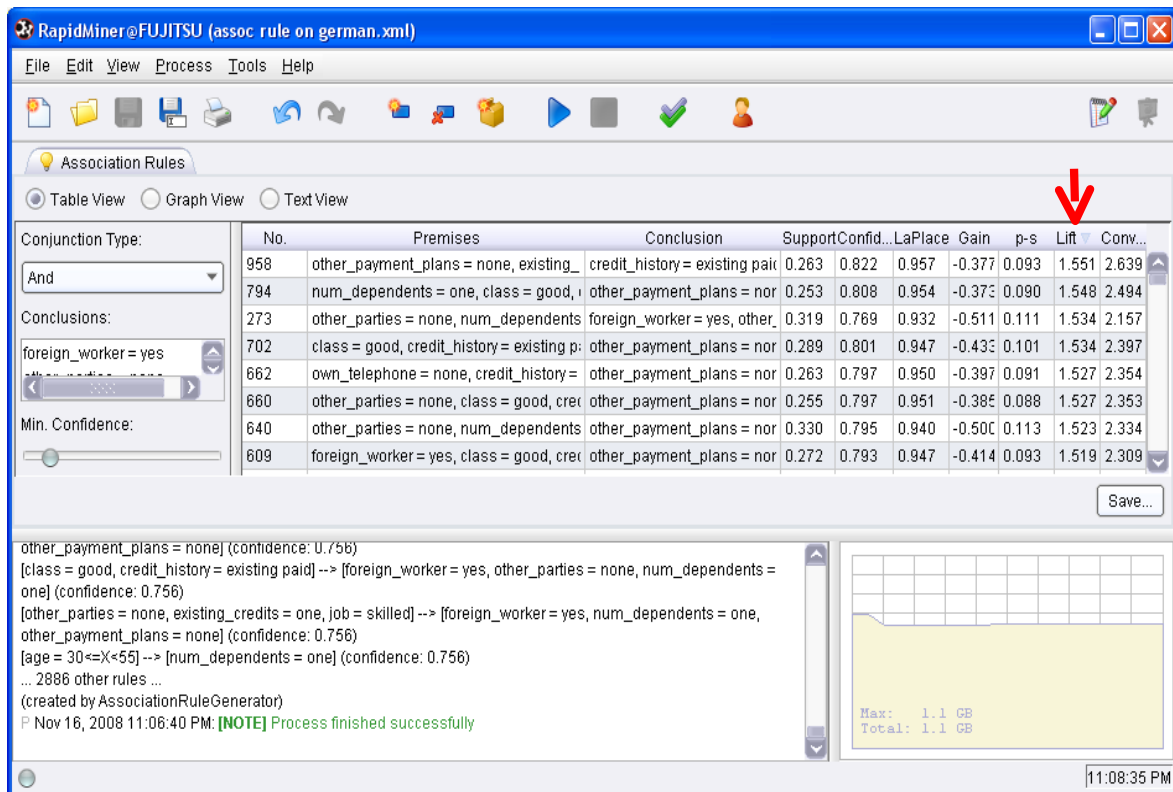
The extraction is carried out in 2 steps. First, with the **FPGROWTH** component, we generate the frequent itemsets. The settings must be defined carefully.



Then, with the **ASSOCIATIONRULEGENERATOR**, we generate the rules from the itemsets.



Launching the computation and visualization of rules. We launch the calculations by clicking on the PLAY menu into the toolbar.



We obtain 2986 rules, the same as the A PRIORI component of Tanagra. Various indicators enable to rank rules. Below, we sort the rules according to the LIFT criterion.

RapidMiner@FUJITSU (assoc rule on german.xml)

File Edit View Process Tools Help

Association Rules

Table View Graph View Text View

Conjunction Type: And

Conclusions:

- foreign_worker = yes
- other_parties = none
- num_dependents = one
- other_payment_plans = none
- housing = own
- class = good
- existing_credits = one
- credit_history = existing paid

Min. Confidence: [Slider]

No.	Premises	Conclusion	Supp...	Conf...	LaPl...	Gain	p-s	Lift	Con...
243	other_parties = none, other_payment_plans = none, checking_status	class = good	0.290	0.927	0.982	-0.33	0.071	1.32	4.08
241	foreign_worker = yes, other_parties = none, other_payment_plans = r	class = good	0.280	0.924	0.982	-0.32	0.068	1.32	3.95
241	other_parties = none, num_dependents = one, other_payment_plans	class = good	0.250	0.923	0.982	-0.29	0.060	1.31	3.87
237	other_payment_plans = none, checking_status = no checking	class = good	0.303	0.916	0.980	-0.35	0.072	1.31	3.66
233	num_dependents = one, other_payment_plans = none, checking_sta	class = good	0.261	0.916	0.981	-0.30	0.062	1.30	3.56
231	foreign_worker = yes, other_payment_plans = none, checking_status	class = good	0.291	0.915	0.980	-0.34	0.068	1.30	3.53
228	foreign_worker = yes, num_dependents = one, other_payment_plans	class = good	0.253	0.913	0.981	-0.30	0.058	1.30	3.46
205	other_parties = none, housing = own, checking_status = no checking	class = good	0.259	0.896	0.977	-0.31	0.057	1.28	2.89
203	housing = own, checking_status = no checking	class = good	0.272	0.896	0.976	-0.33	0.056	1.27	2.85
201	foreign_worker = yes, other_parties = none, housing = own, checking	class = good	0.251	0.893	0.977	-0.31	0.054	1.27	2.81
198	foreign_worker = yes, housing = own, checking_status = no checking	class = good	0.263	0.892	0.976	-0.32	0.057	1.27	2.76

Save...

other_payment_plans = none] (confidence: 0.756)
[class = good, credit_history = existing paid] --> [foreign_worker = yes, other_parties = none, num_dependents = one] (confidence: 0.756)
[other_parties = none, existing_credits = one, job = skilled] --> [foreign_worker = yes, num_dependents = one, other_payment_plans = none] (confidence: 0.756)
[age = 30<=X<55] --> [num_dependents = one] (confidence: 0.756)
... 2886 other rules ...
(created by AssociationRuleGenerator)
P Nov 16, 2008 11:06:40 PM: [NOTE] Process finished successfully

Max: 1.1 GB
Total: 1.1 GB

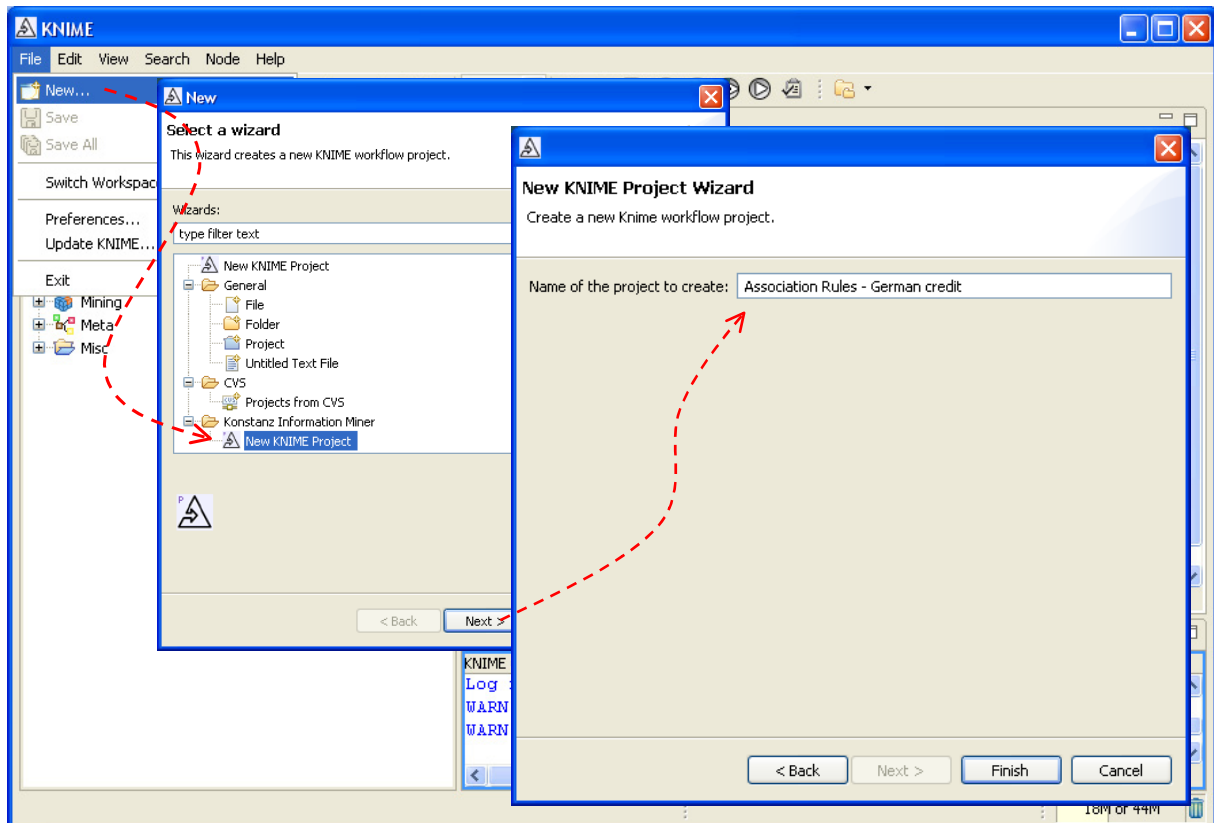
11:16:47 PM

Another option is available. We can filter out the rules according to the presence or absence of an item or a set of items. This is very useful.

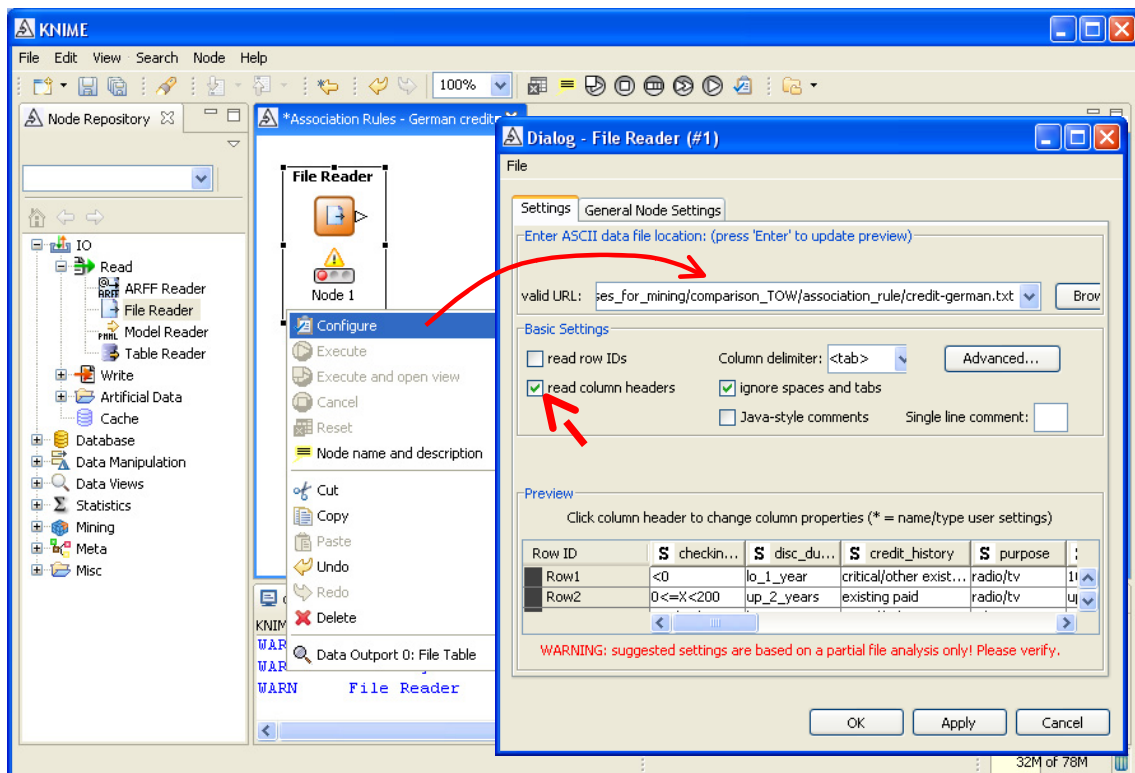
8 Knime

A double data preparation is necessary for KNIME before launching the learning algorithm: a coding o/1 of the attribute-value dataset, followed by a transformation into transactional data.

Workflow creation and data importation. We create a new project by clicking on the FILE / NEW menu. We choose a KNIME project and we specify the name of the project.

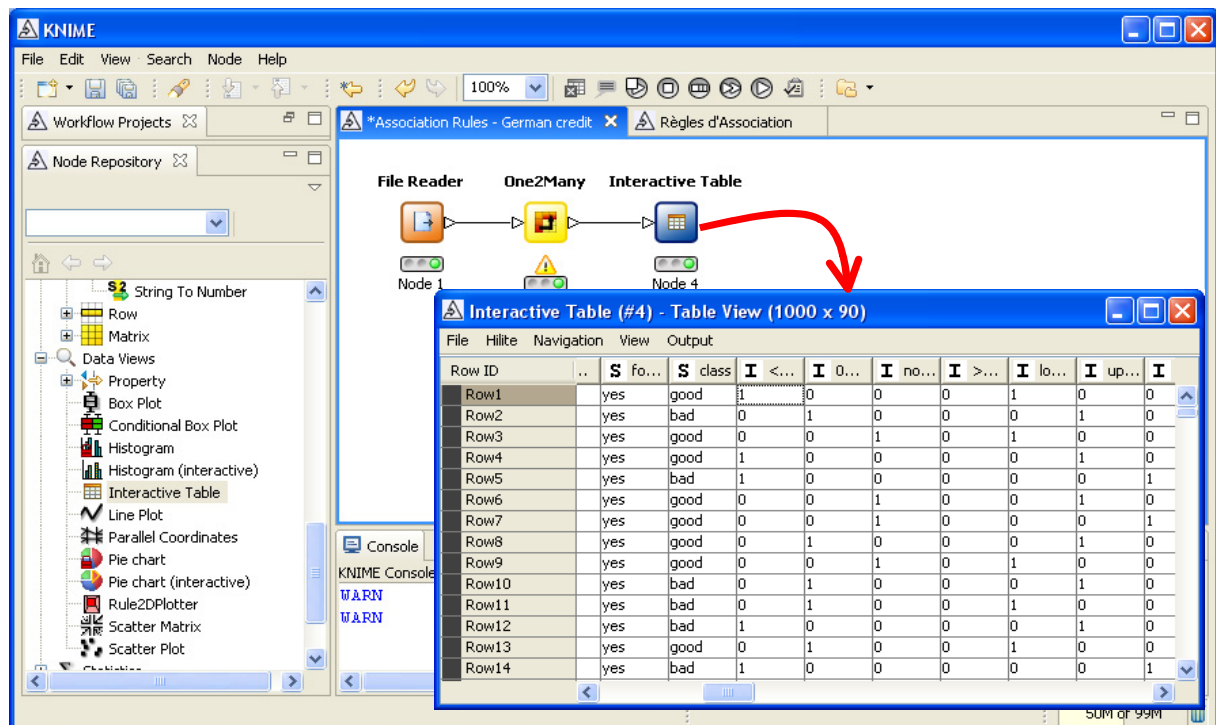


An empty workflow appears. We insert the file access component (FILE READER). We select our data file (CONFIGURE menu). We click on the EXECUTE menu in order to load the dataset.



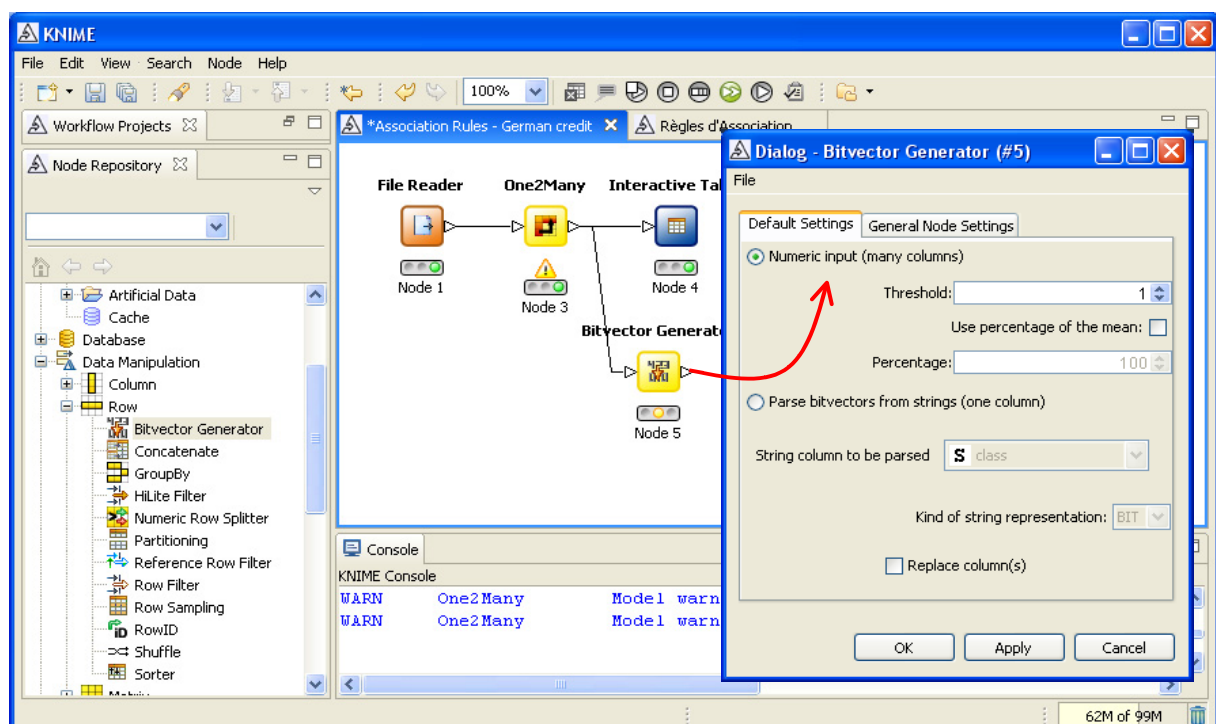
First coding step. First we must recode the attribute-value dataset into a binary dataset. We use the ONE2MANY component (DATA MANIPULATION / COLUMN). We connect this component to the

previous one. We add an INTERACTIVE TABLE component (DATA VIEWS), that we connect to ONE2MANY, in order to view the new dataset.

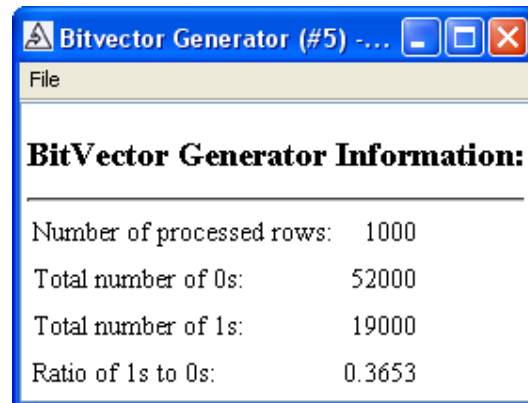


There are 90 columns now. To the previous 19 variables are added 71 binary variables.

Second coding step. The first transformation is not enough. We must go through an internal format specific using the BITVECTOR component (DATA MANIPULATION / COLUMN). We connect it to ONE2MANY. We click on the CONFIGURE menu, we select only the binary variables.



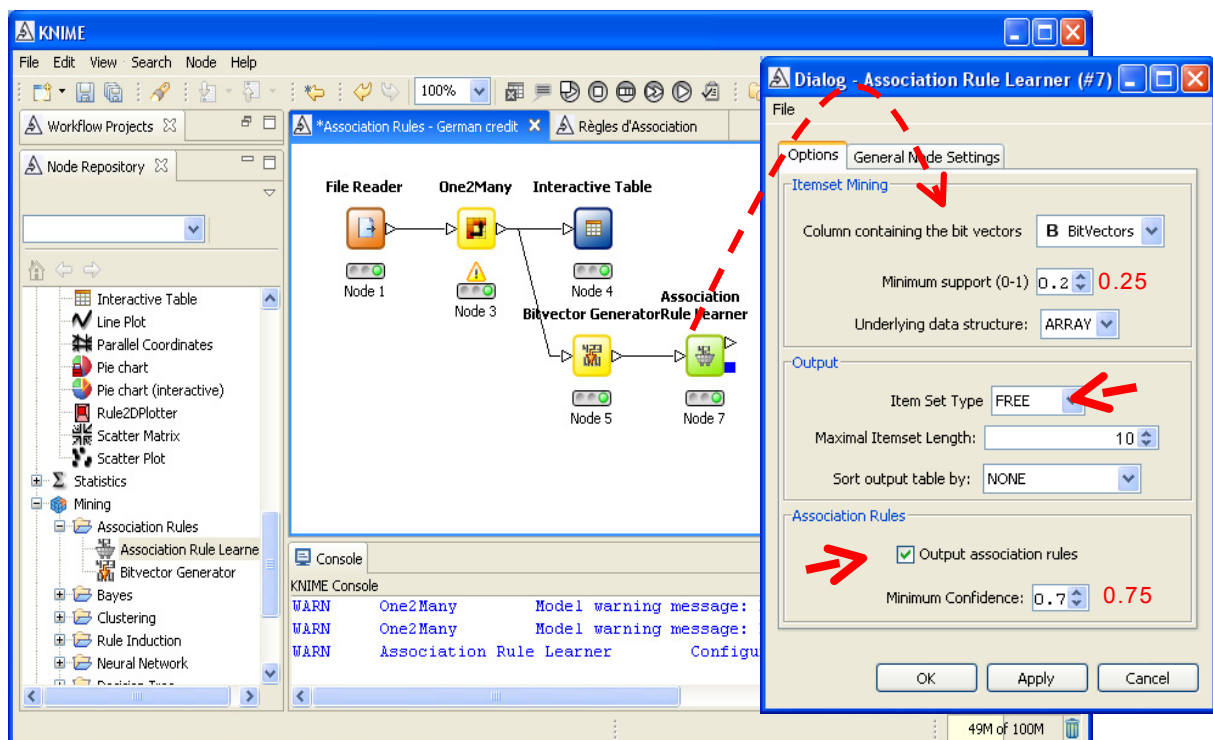
By clicking on the EXECUTE AND OPEN VIEW menu, we obtain a description of the generated transactional dataset.



BitVector Generator Information:	
Number of processed rows:	1000
Total number of 0s:	52000
Total number of 1s:	19000
Ratio of 1s to 0s:	0.3653

To obtain the information of density of R, we make $19000/(19000+52000) = 26.76\%$.

Extraction of rules. We use the ASSOCIATION RULE LEARNER component in order to extract rules. We set the following parameters.

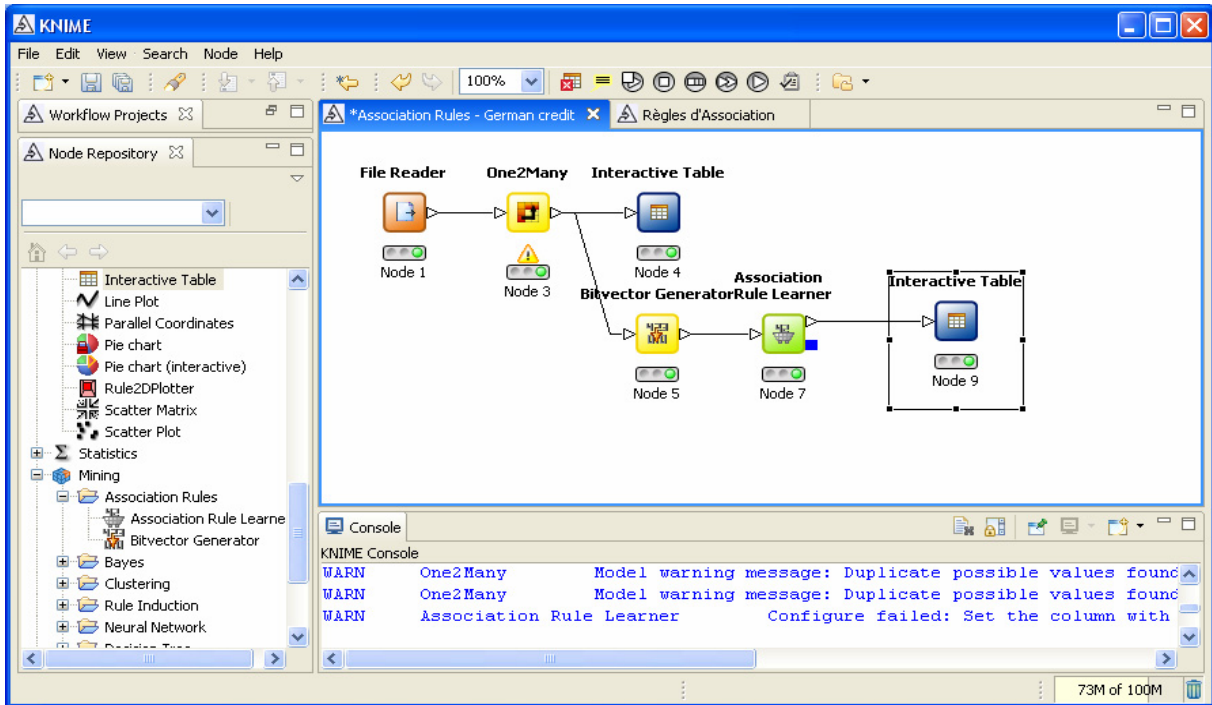


The screenshot shows the KNIME software interface with a workflow and the Association Rule Learner dialog box open. The workflow consists of the following nodes: File Reader (Node 1), One2Many (Node 3), Interactive Table (Node 4), Bitvector Generator (Node 5), and Association Rule Learner (Node 7). The Association Rule Learner dialog box is configured with the following parameters:

- Options: General Mode Settings
- Itemset Mining:
 - Column containing the bit vectors: B BitVectors
 - Minimum support (0-1): 0.2 (set to 0.25)
 - Underlying data structure: ARRAY
- Output:
 - Item Set Type: FREE
 - Maximal Itemset Length: 10
 - Sort output table by: NONE
- Association Rules:
 - Output association rules
 - Minimum Confidence: 0.7 (set to 0.75)

Red arrows in the image point to the 'Item Set Type' dropdown, the 'Output association rules' checkbox, and the 'Minimum Confidence' field.

The INTERACTIVE TABLE component enables to visualize the rules.



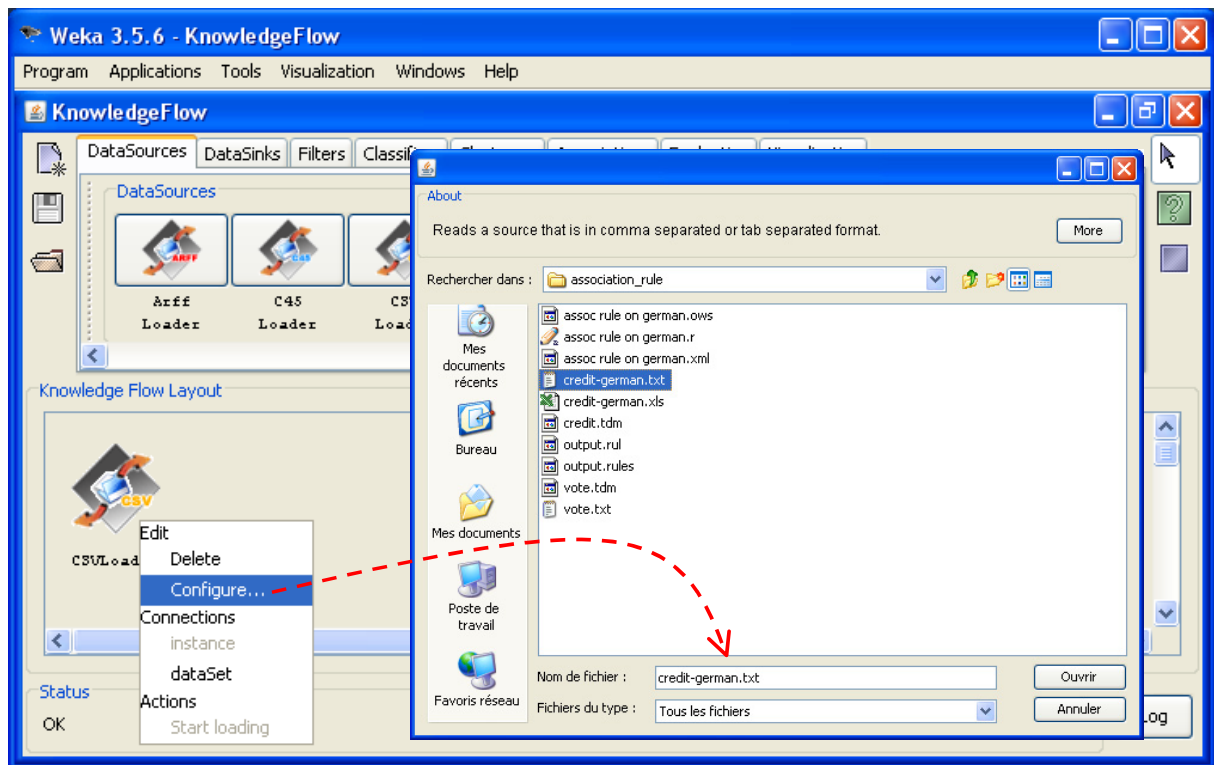
KNIME extracts only the rules with one item into the consequent. 1928 rules are generated. The visualization tool has not interactive functionalities (ranking or filtering rules).

Row ID	D Support	D Confide...	S Consequent	S Implies	S Item0	S Item1	S Item2	S Item3
rule0	0.259	0.945	yes_foreign_worker	<---	<0_checking_sta...	?	?	?
rule1	0.264	0.981	yes_foreign_worker	<---	0<=X<200_chec...	?	?	?
rule2	0.373	0.947	none_other_parties	<---	no_checking_chec...	?	?	?
rule3	0.313	0.948	none_other_parties	<---	no_checking_chec...	none_other...	?	?
rule4	0.313	0.839	none_other_payment_plans	<---	no_checking_chec...	none_other...	?	?
rule5	0.271	0.951	none_other_parties	<---	no_checking_chec...	none_other...	one_num_d...	?
rule6	0.271	0.86	none_other_payment_plans	<---	no_checking_chec...	none_other...	one_num_d...	?
rule7	0.271	0.866	one_num_dependents	<---	no_checking_chec...	none_other...	none_other...	?
rule8	0.264	0.953	none_other_parties	<---	no_checking_chec...	none_other...	one_num_d...	yes_foreign...
rule9	0.264	0.86	none_other_payment_plans	<---	no_checking_chec...	none_other...	one_num_d...	yes_foreign...
rule10	0.264	0.871	one_num_dependents	<---	no_checking_chec...	none_other...	none_other...	yes_foreign...
rule11	0.264	0.974	yes_foreign_worker	<---	no_checking_chec...	none_other...	none_other...	one_num_d...
rule12	0.25	0.958	none_other_parties	<---	no_checking_chec...	none_other...	one_num_d...	good_class
rule13	0.25	0.893	none_other_payment_plans	<---	no_checking_chec...	none_other...	yes_foreign...	good_class
rule14	0.25	0.862	one_num_dependents	<---	no_checking_chec...	none_other...	none_other...	good_class
rule15	0.25	0.923	good_class	<---	no_checking_chec...	none_other...	none_other...	one_num_d...
rule16	0.303	0.953	none_other_parties	<---	no_checking_chec...	none_other...	yes_foreign...	?
rule17	0.303	0.837	none_other_payment_plans	<---	no_checking_chec...	none_other...	yes_foreign...	?
rule18	0.303	0.968	yes_foreign_worker	<---	no_checking_chec...	none_other...	none_other...	?
rule19	0.28	0.962	none_other_parties	<---	no_checking_chec...	none_other...	yes_foreign...	good_class
rule20	0.28	0.875	none_other_payment_plans	<---	no_checking_chec...	none_other...	yes_foreign...	good_class
rule21	0.28	0.966	yes_foreign_worker	<---	no_checking_chec...	none_other...	none_other...	good_class
rule22	0.28	0.924	good_class	<---	no_checking_chec...	none_other...	none_other...	yes_foreign...
rule23	0.28	0.957	none_other_parties	<---	no_checking_chec...	none_other...	good_class	?

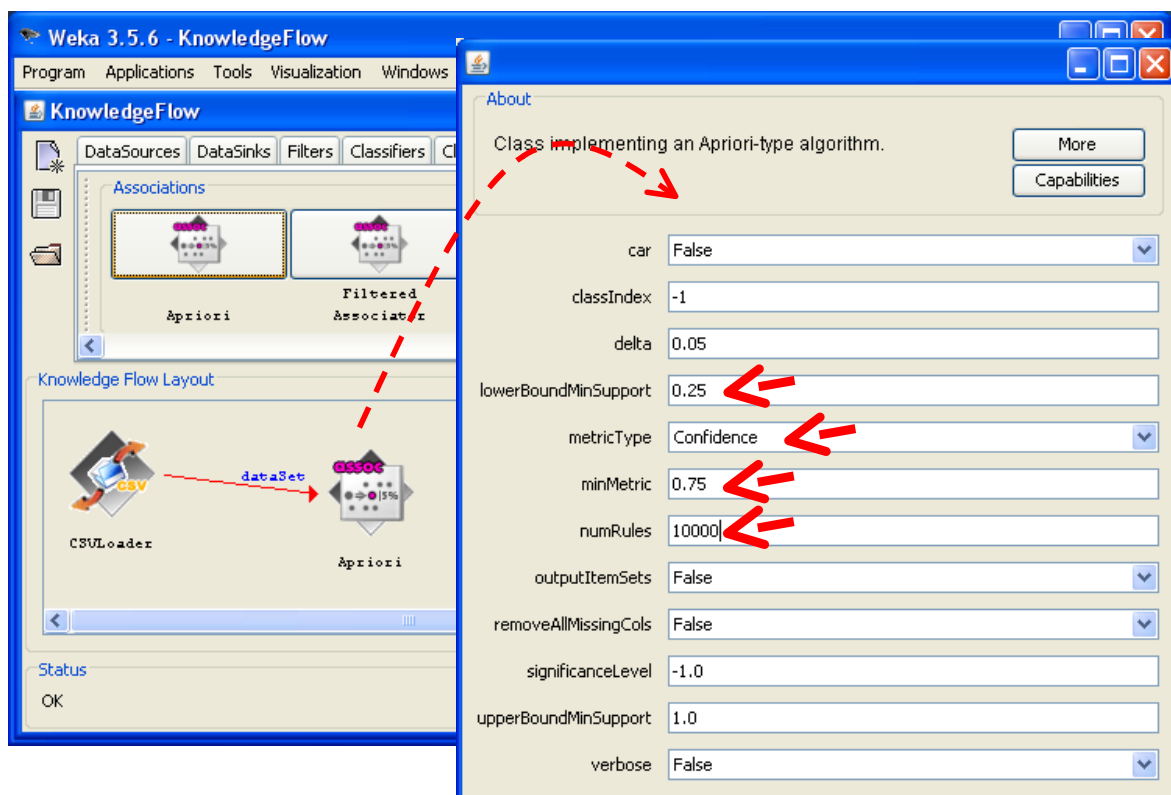
9 Weka

Like RapidMiner, it is more convenient to define the whole treatment before launching the computation with Weka. We use the Knowledge Flow in this tutorial.

Specifying the operations. When we launch Weka, we click on the APPLICATIONS / KNOWLEDGEFLOW menu. First, we must insert the data access component CSV LOADER (DATASOURCES tab). We select the data file.

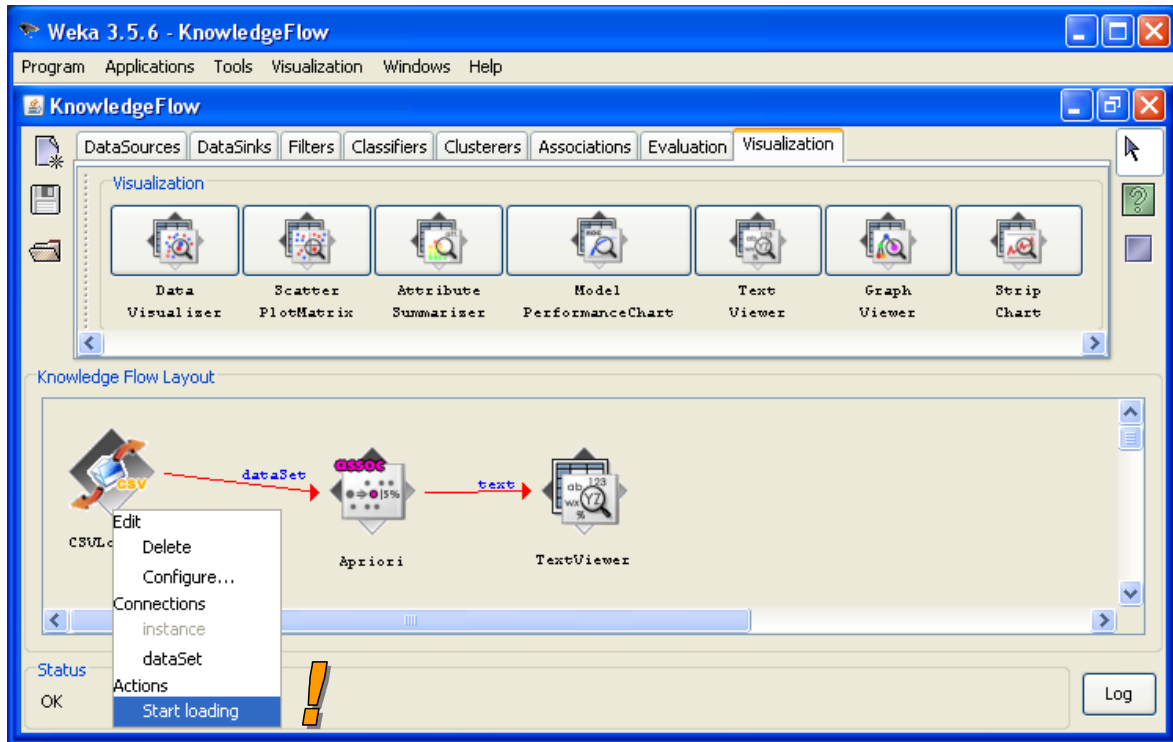


We add the APRIORI component (ASSOCIATIONS tab). We connect the previous component to this one. We click on the CONFIGURE menu in order to define the parameters.

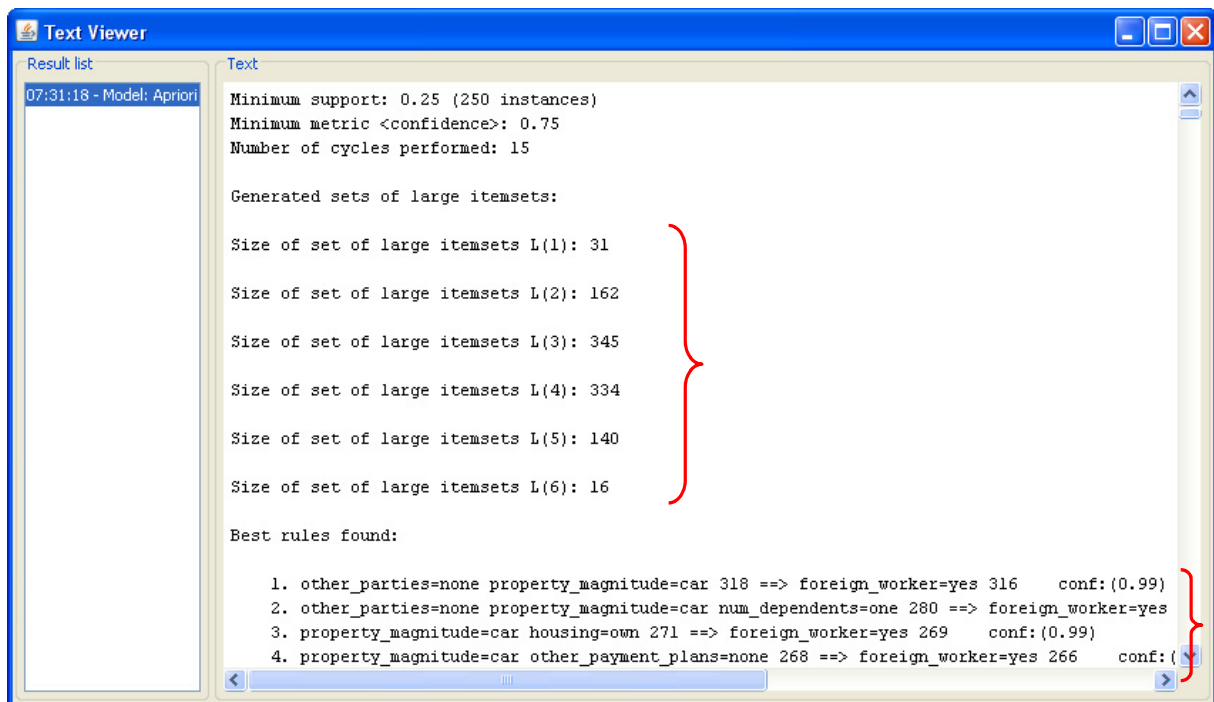


Setting NUMRULES to 1000 removes the restriction delta on the number of generated rules. We add then the TEXT VIEWER component in order to visualize the rules.

We launch the calculations by clicking on the START LOADING menu of... the CSV LOADER component into the knowledge flow.



To visualize the results, we click on the SHOW RESULTS menu of the TEXTVIEWER component. Like the A PRIORI component of Tanagra, we obtain 2986 rules.



10 Conclusion

All software presented in this tutorial can extract association rules from a data file in an "individuals x variables" format. For some of them, a data preparation is required before to produce a transactional data format. This is not always obvious, especially when the software is not well documented. I admit to having groping a bit. But finally, once the data are correctly generated, the software produced similar results. This is what matters.