

1 Theme

Data Mining with R – The “rattle” package.

R (<http://www.r-project.org/>) is one of the most exciting free data mining software projects of these last years. Its popularity is completely justified (see [Kdnuggets Polls – Data Mining/ Analytic Tools Used – 2011](#)). Among the reasons which explain this success, we distinguish two very interesting characteristics: (1) we can extend almost indefinitely the features of the tool with the packages; (2) we have a programming language which allows to perform easily sequences of complex operations.

But this second property can be also a drawback. Indeed, some users do not want to learn a new programming language before being able to realize projects. For this reason, tools which allow to define the sequence of commands with diagrams (such as Tanagra, Knime, RapidMiner, etc.) still remain a valuable alternative with the data miners.

In this tutorial, we present the "Rattle" package which allows to the data miners to use R without needing to know the associated programming language. All the operations are performed with simple clicks, such as for any software driven by menus. But, in addition, all the commands are stored. We can save them in a file. Then, in a new working session, we can easily repeat all the operations. Thus, we find one of the important properties which miss to the tools driven by menus.

To describe the use of the rattle package, we perform an analysis similar to the one suggested by the rattle's author in its presentation paper (G.J. Williams, « Rattle : A Data Mining GUI for R », in *The R Journal*, volume 1 / 2, pages 45—55, December 2009, http://journal.r-project.org/archive/2009-2/RJournal_2009-2_Williams.pdf). We perform the following steps: loading the data file; partitioning the instances into learning and test samples; specifying the types of the variables (target or input); computing some descriptive statistics; learning the predictive models from the learning sample; assessing the models on the test sample (confusion matrix, error rate, some curves).

2 Dataset

We use the «heart»¹ data file. We want to explain the occurrence of the DISEASE from the characteristics of patients. We show here the first instances of the dataset.

	A	B	C	D	E	F	G	H	I	
1	age	sex	chest_pain	trestbps	chol	fbs	restecg	thalach	exang	disease
2	31	male	asympt	120	270	f	normal	153	yes	positive
3	33	female	asympt	100	246	f	normal	150	yes	positive
4	34	male	typ_angina	140	156	f	normal	180	no	positive
5	35	male	atyp_angina	110	257	f	normal	140	no	positive
6	36	male	atyp_angina	120	267	f	normal	160	no	positive
7	37	male	asympt	140	207	f	normal	130	yes	positive
8	38	male	asympt	110	196	f	normal	166	no	positive
9	38	male	asympt	120	282	f	normal	170	no	positive
10	38	male	asympt	92	117	f	normal	134	yes	positive
11	41	male	asympt	110	289	f	normal	170	no	positive
12	43	male	asympt	150	247	f	normal	130	yes	positive

¹ http://eric.univ-lyon2.fr/~ricco/tanagra/fichiers/heart_for_rattle.txt ; a description of this data file is available on the following website: <http://archive.ics.uci.edu/ml/datasets/Heart+Disease>

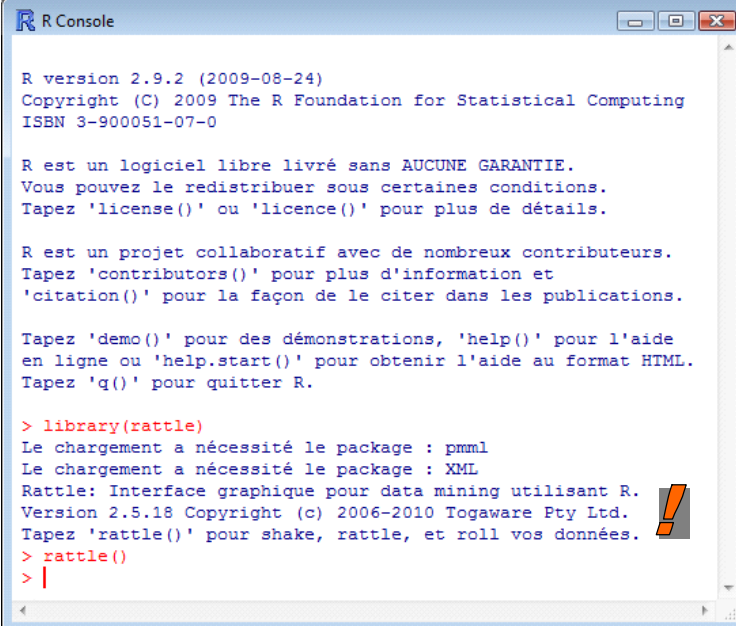
3 Data Mining with Rattle

3.1 Loading the rattle package

First, we load the rattle package [`library()`]. Then, we start the GUI with the command `rattle()`.

```
> #loading the package
> library(rattle)
> #lauching the GUI
> rattle()
```

Into the R console, we have...



```
R Console

R version 2.9.2 (2009-08-24)
Copyright (C) 2009 The R Foundation for Statistical Computing
ISBN 3-900051-07-0

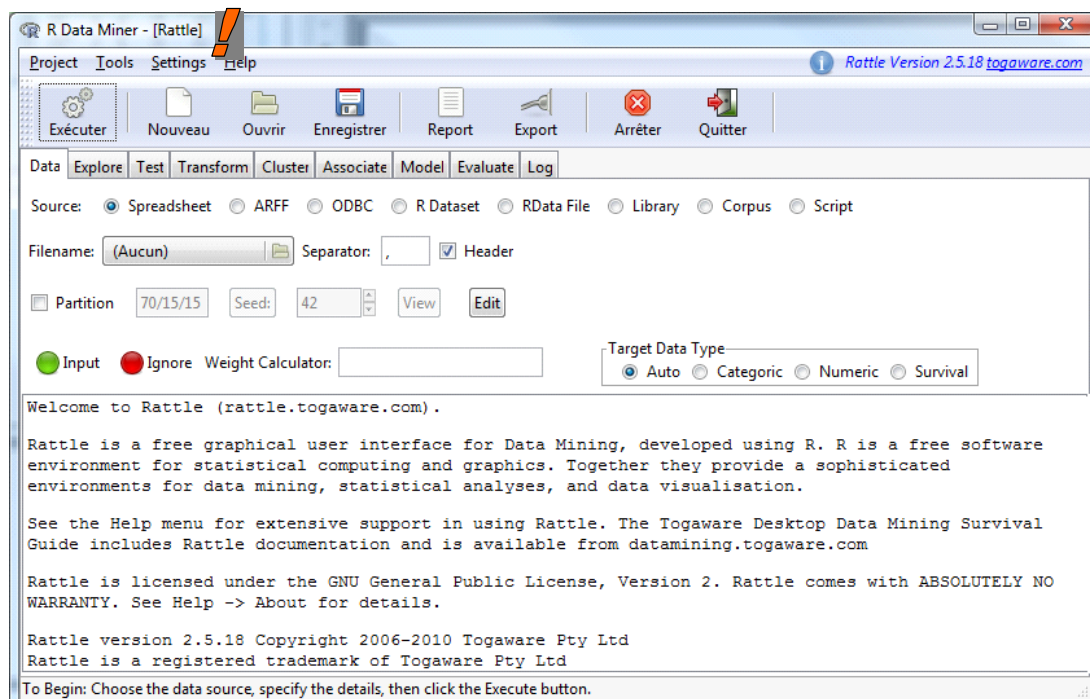
R est un logiciel libre livré sans AUCUNE GARANTIE.
Vous pouvez le redistribuer sous certaines conditions.
Tapez 'license()' ou 'licence()' pour plus de détails.

R est un projet collaboratif avec de nombreux contributeurs.
Tapez 'contributors()' pour plus d'information et
'citation()' pour la façon de le citer dans les publications.

Tapez 'demo()' pour des démonstrations, 'help()' pour l'aide
en ligne ou 'help.start()' pour obtenir l'aide au format HTML.
Tapez 'q()' pour quitter R.

> library(rattle)
Le chargement a nécessité le package : pmml
Le chargement a nécessité le package : XML
Rattle: Interface graphique pour data mining utilisant R.
Version 2.5.18 Copyright (c) 2006-2010 Togaware Pty Ltd.
Tapez 'rattle()' pour shake, rattle, et roll vos données.
> rattle()
> |
```

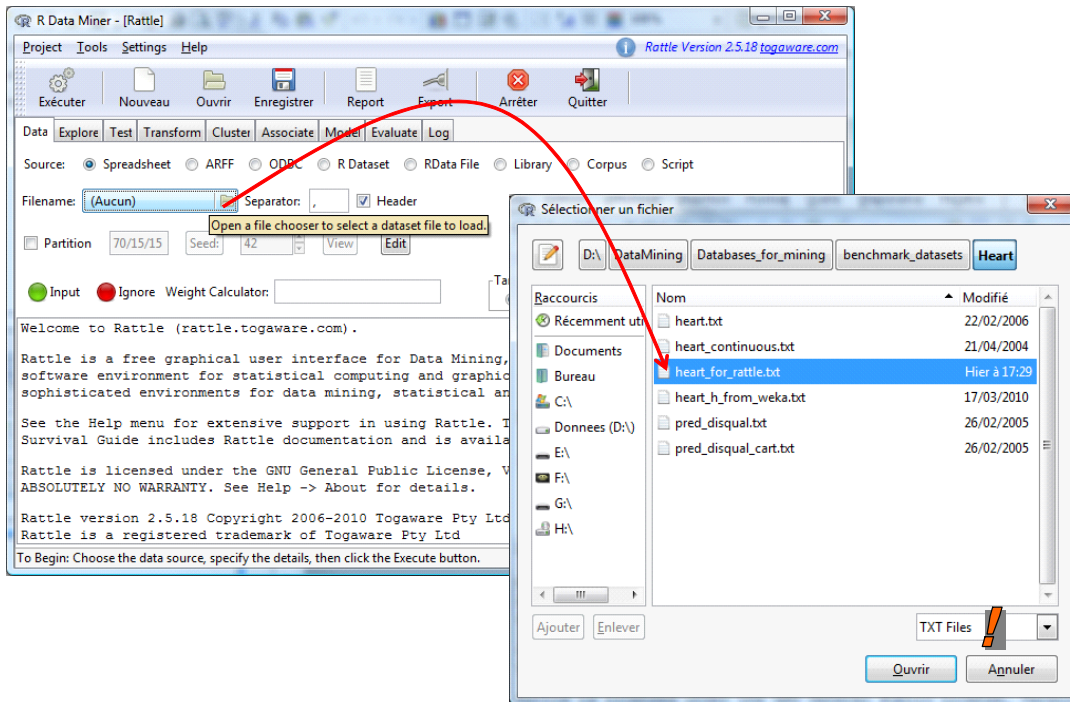
From now, we perform all the operations by clicking on the appropriate menu or button. All these operations are recorded as R commands by rattle. The rattle GUI is displayed.



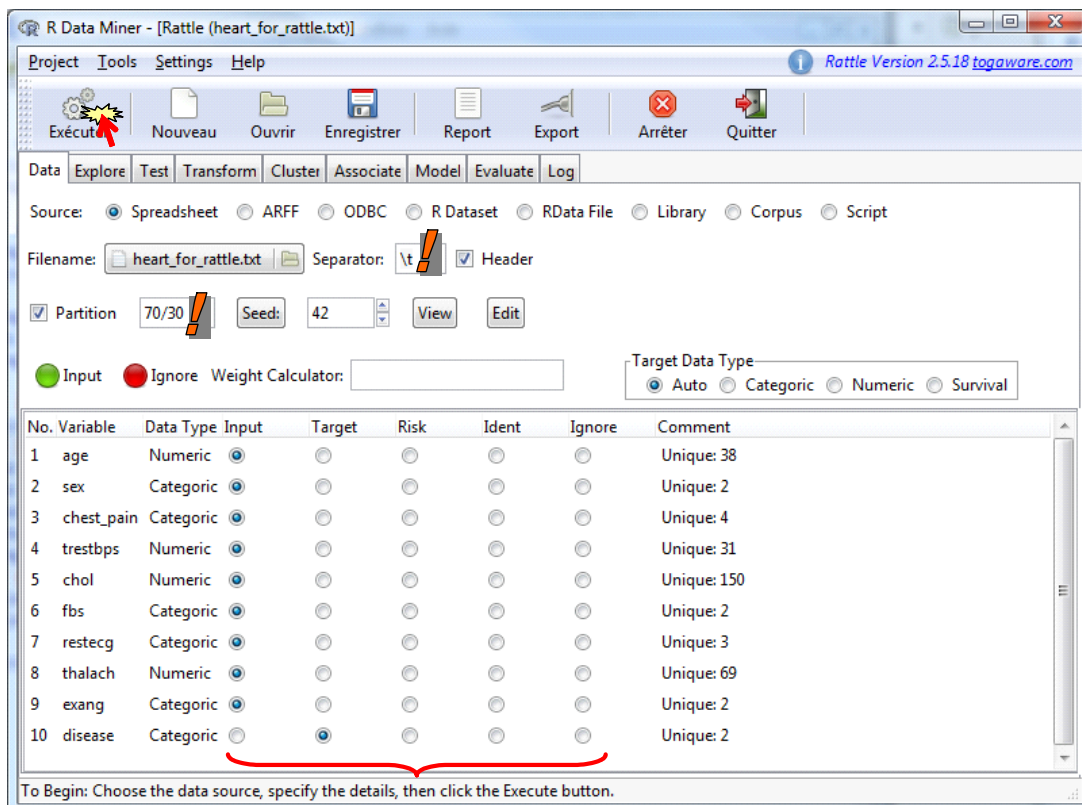
The use of rattle is always the same: we define the command by working in the appropriate tab (Data: load the dataset; Explore: some descriptive statistics; Test: some statistical tests, etc.); then, we launch the calculations by clicking on the EXECUTER button into the toolbar.

3.2 Importing the data file

Into the "Data" tab, we click on the FILENAME button. We select the "heart_for_rattle.txt" data file.



We specify the column separator: « SEPARATOR = \t ». Then we click on EXECUTER.



The dataset is loaded. The variable type is automatically detected from the distinct values into each column (discrete or continuous). We can define the TARGET attribute and the INPUT ones. Last, we specify the size of the training (70% of instances, drawn randomly) and test (30%) samples.

3.3 Dataset description

R Data Miner - [Rattle (heart_for_rattle.txt)]
Rattle Version 2.5.18 togaware.com

Project Tools Settings Help

Exécuter Nouveau Ouvrir Enregistrer Report Export Arrêter Quitter

Data Explore Transform Cluster Associate Model Evaluate Log

Type: Summary Distributions Correlation Principal Components Interactive

Summary Describe Basics Kurtosis Skewness Show Missing

Below is a summary of the dataset.
The data is limited to the training dataset.
Data frame: crs\$dataset[crs\$sample,] 200 observations and 10 variables Maximum # NAs:0

Variable	Levels	Storage
age		integer
sex	2	integer
chest_pain	4	integer
trestbps		integer
chol		integer
fbs	2	integer
restecg	3	integer
thalach		integer
exang	2	integer
disease	2	integer

```

+-----+
|Variable|Levels|
+-----+
|sex     |female,male|
+-----+
|chest_pain|asympt,atyp_angina,non_anginal,typ_angina|
+-----+
|fbs     |f,t|
+-----+
|restecg |left_vent_hyper,normal,st_t_wave_abnormality|
+-----+
|exang   |no,yes|
+-----+
|disease |negative,positive|
+-----+

```

For the simple distribution tables below the 1st and 3rd Qu. refer to the first and third quartiles, indicating that 25% of the observations have values of or greater than (respectively) the value listed.

age	sex	chest_pain	trestbps	chol
Min. :28.00	female: 53	asympt :81	Min. : 98.0	Min. :132.0
1st Qu.:42.00	male :147	atyp_angina:77	1st Qu.:120.0	1st Qu.:211.0
Median :49.00		non_anginal:34	Median :130.0	Median :250.0
Mean :48.27		typ_angina : 8	Mean :133.7	Mean :252.0
3rd Qu.:54.00			3rd Qu.:140.0	3rd Qu.:277.5
Max. :65.00			Max. :190.0	Max. :603.0

fbs	restecg	thalach	exang	disease
f:184	left_vent_hyper : 5	Min. : 82.0	no :142	negative:135
t: 16	normal :157	1st Qu.:122.0	yes: 58	positive: 65
	st_t_wave_abnormality: 38	Median :140.0		
		Mean :139.6		
		3rd Qu.:155.2		
		Max. :190.0		

Generated by Rattle 2010-06-15 09:50:20 Maison

Find: Rechercher Suivant

Data summary generated.

Into the Explore tab, we obtain some descriptive statistics indicators about the variables (SUMMARY / SUMMARY option). For the discrete variables, rattle lists the values (levels). For the continuous ones, we have the min, max, mean, quartiles. **All the indicators are computed on the learning sample.**

With the SUMMARY / DESCRIBE option, we obtain a more detailed description. Among others, for the continuous variables, the indications are useful to detect unusual values (outliers).

R Data Miner - [Rattle (heart_for_rattle.txt)]

Project Tools Settings Help

Rattle Version 2.5.18 togaware.com

Exécuter Nouveau Ouvrir Enregistrer Report Export Arrêter Quitter

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: Summary Distributions Correlation Principal Components Interactive

Summary Describe Basics Kurtosis Skewness Show Missing

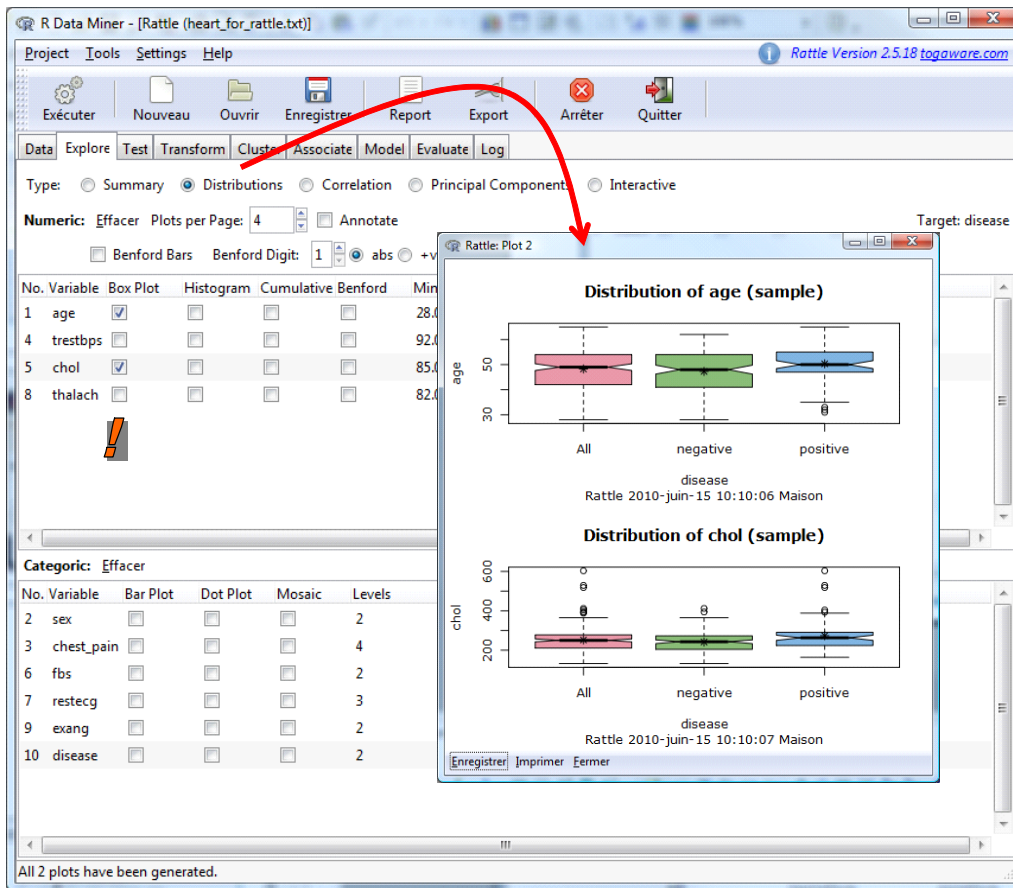
Below is a description of the dataset.
The data is limited to the training dataset.

```
crs$dataset[crs$sample, ]
  10 Variables      200 Observations
-----
age
  n missing unique   Mean   .05   .10   .25   .50   .75   .90   .95
  200      0     37  48.27  34.00  37.00  42.00  49.00  54.00  58.00  59.05
lowest : 28 29 30 31 32, highest: 60 61 62 63 65
-----
sex
  n missing unique
  200      0      2
female (53, 26%), male (147, 74%)
-----
chest_pain
  n missing unique
  200      0      4
asympt (81, 40%), atyp_angina (77, 38%), non_anginal (34, 17%), typ_angina (8, 4%)
-----
trestbps
  n missing unique   Mean   .05   .10   .25   .50   .75   .90   .95
  200      0     25  133.7  110.0  110.0  120.0  130.0  140.0  160.0  160.5
lowest :  98 100 105 106 108, highest: 150 160 170 180 190
-----
```

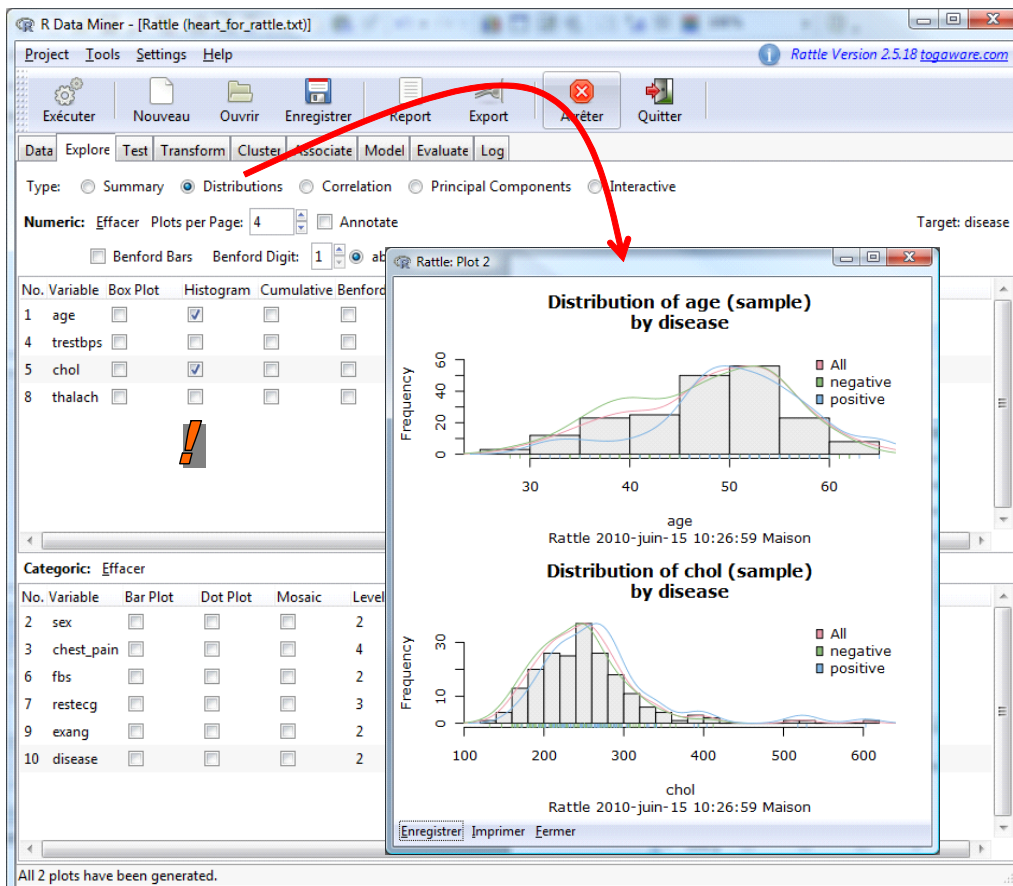
Find: Rechercher Suivant

Data summary generated.

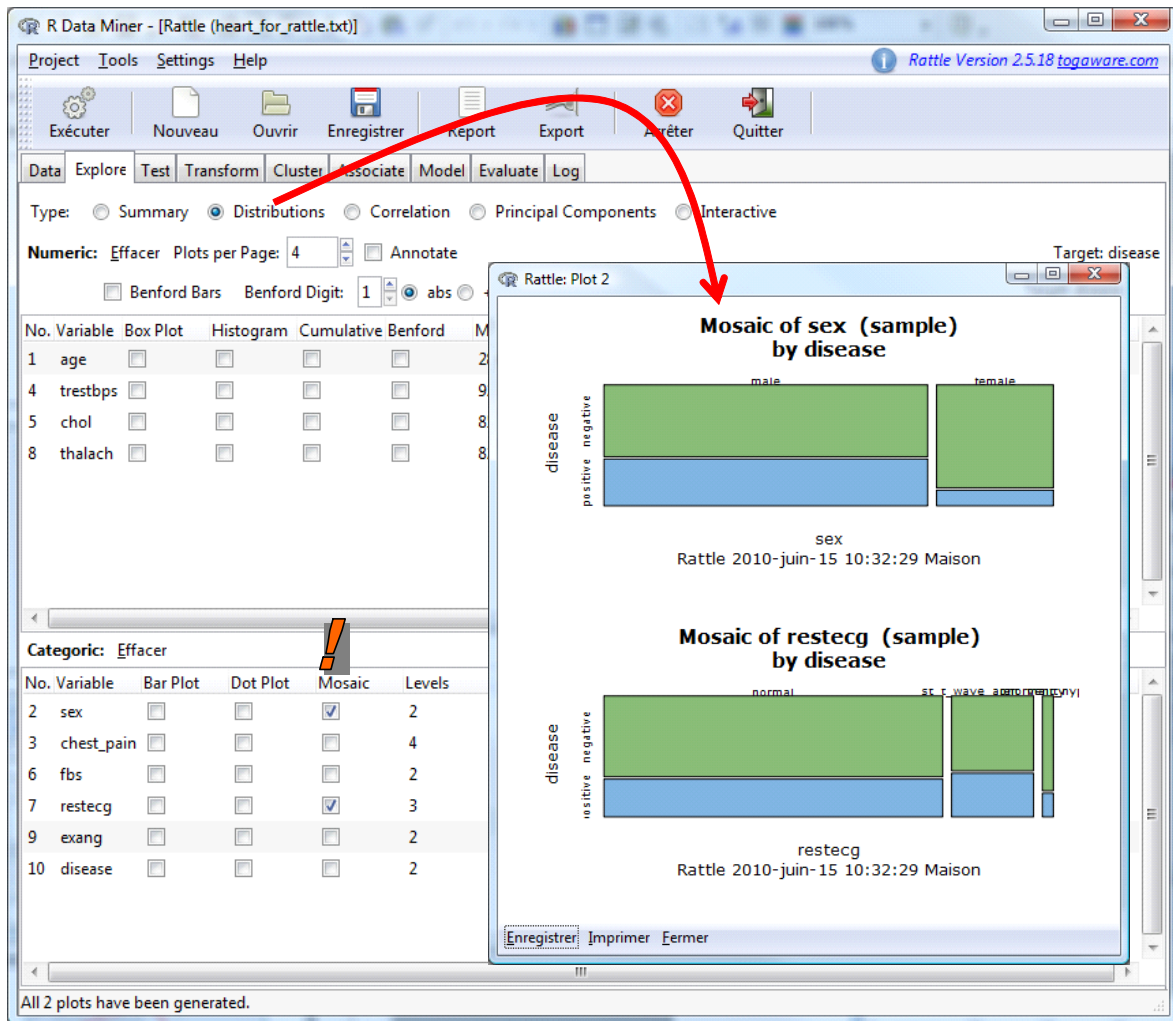
Into the Explore tab still, with the DISTRIBUTIONS option, we obtain some graphical representations of the distributions. We have for instance the conditional box plots of AGE and CHOL according to the values of DISEASE.



We can obtain also the conditional distribution functions.

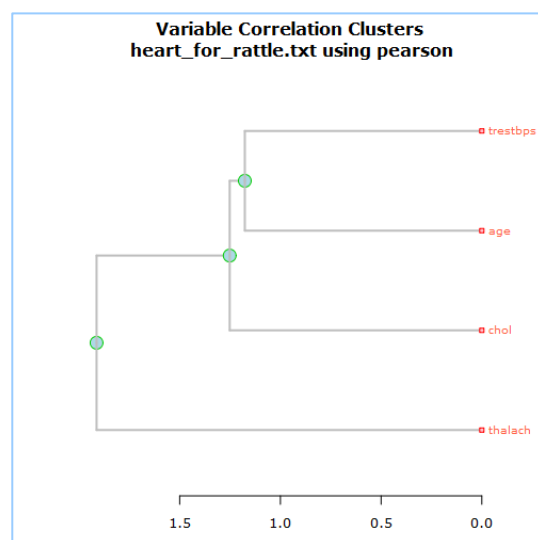


About the discrete variables, we can obtain the "Mosaic" of the variables, according still to the values of the target attribute.



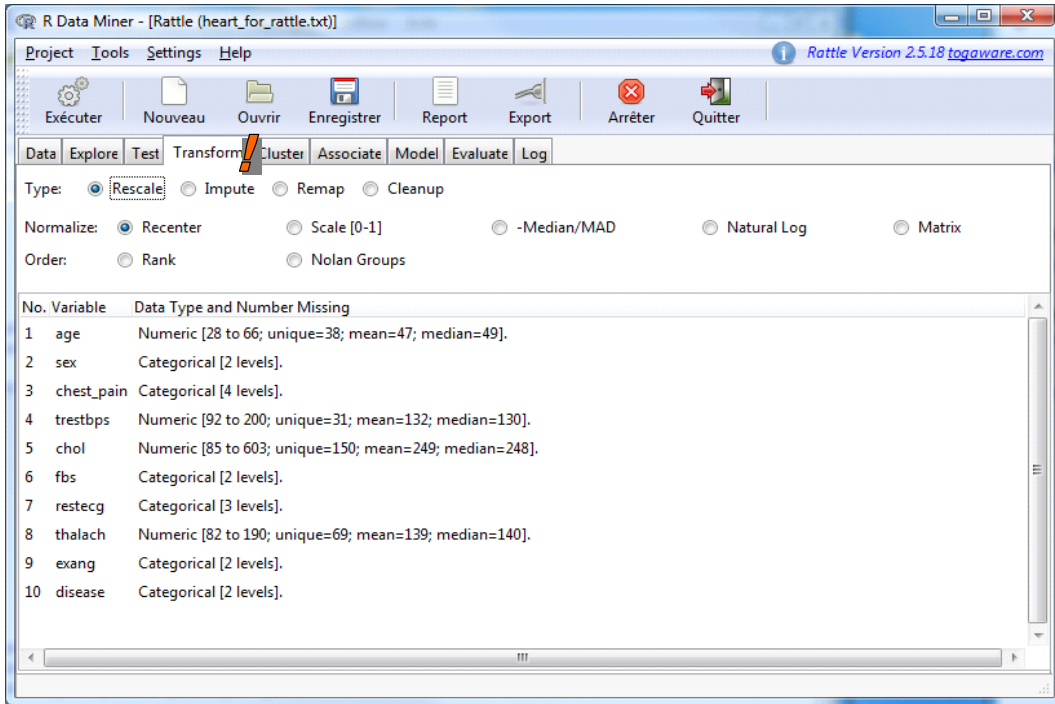
For instance, about SEX, the men (MALE) are more numerous than women (FEMALE) into the sample; and the proportion of disease is higher for the men.

We can also obtain the correlations about the continuous input attributes. The correlations are described in a hierarchical structure. It is useful for instance for the detection of the redundant variables.



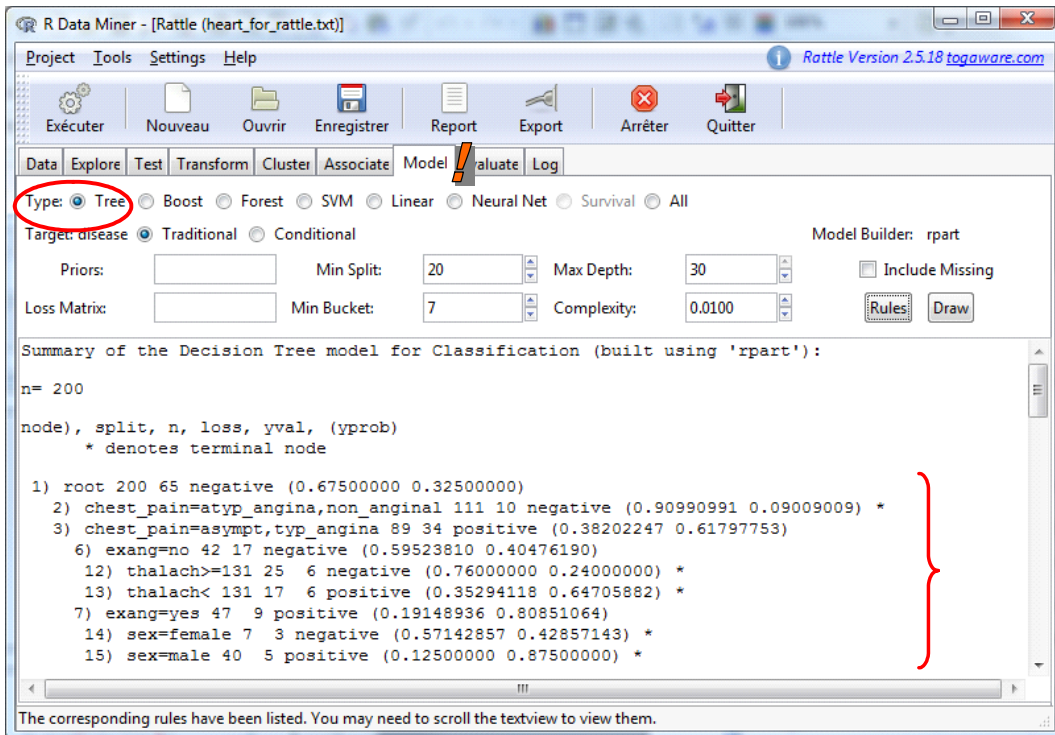
3.4 Data transformation

The "Transform" tab is dedicated to the variable transformation. Some usual operators are available (e.g. logarithm, rank, etc.).



3.5 Supervised learning

This step is at the heart of our analysis. We select the "Model" tab. We want to evaluate three methods: decision tree induction, random forest, logistic regression.



About the **decision tree**, rattle uses the rpart command from the rpart package. We note the default parameters used. We click on the EXECUTER button. We obtain the rules associated to the tree by clicking on the RULES button.

Tree as rules:

```

Rule number: 15 [yval=positive cover=40 (20%) prob=0.88]
chest_pain=asympt,typ_angina
exang=yes
sex=male

Rule number: 13 [yval=positive cover=17 (8%) prob=0.65]
chest_pain=asympt,typ_angina
exang=no
thalach< 131

Rule number: 14 [yval=negative cover=7 (4%) prob=0.43]
chest_pain=asympt,typ_angina
exang=yes
sex=female

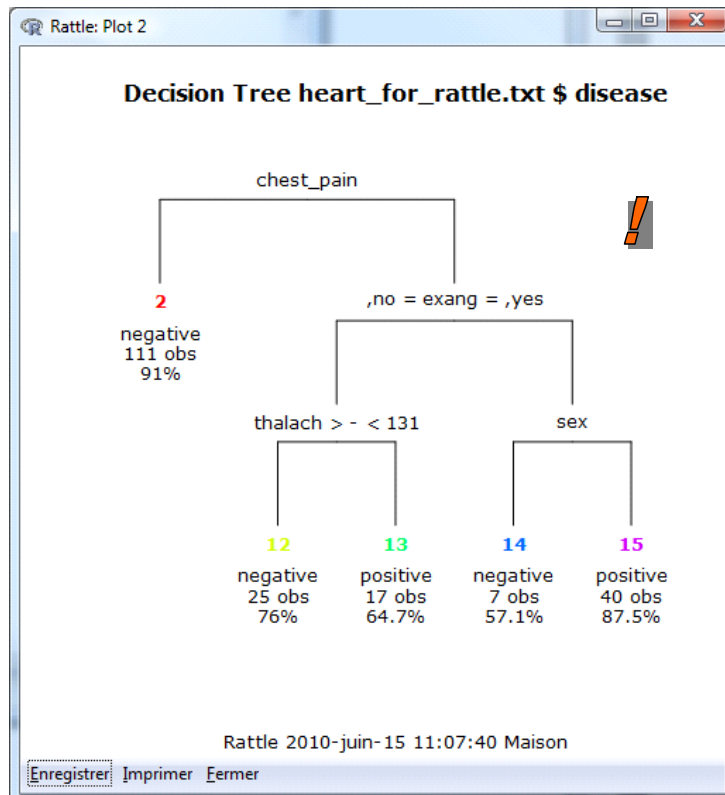
Rule number: 12 [yval=negative cover=25 (12%) prob=0.24]
chest_pain=asympt,typ_angina
exang=no
thalach>=131

Rule number: 2 [yval=negative cover=111 (56%) prob=0.09]
chest_pain=atyp_angina,non_anginal

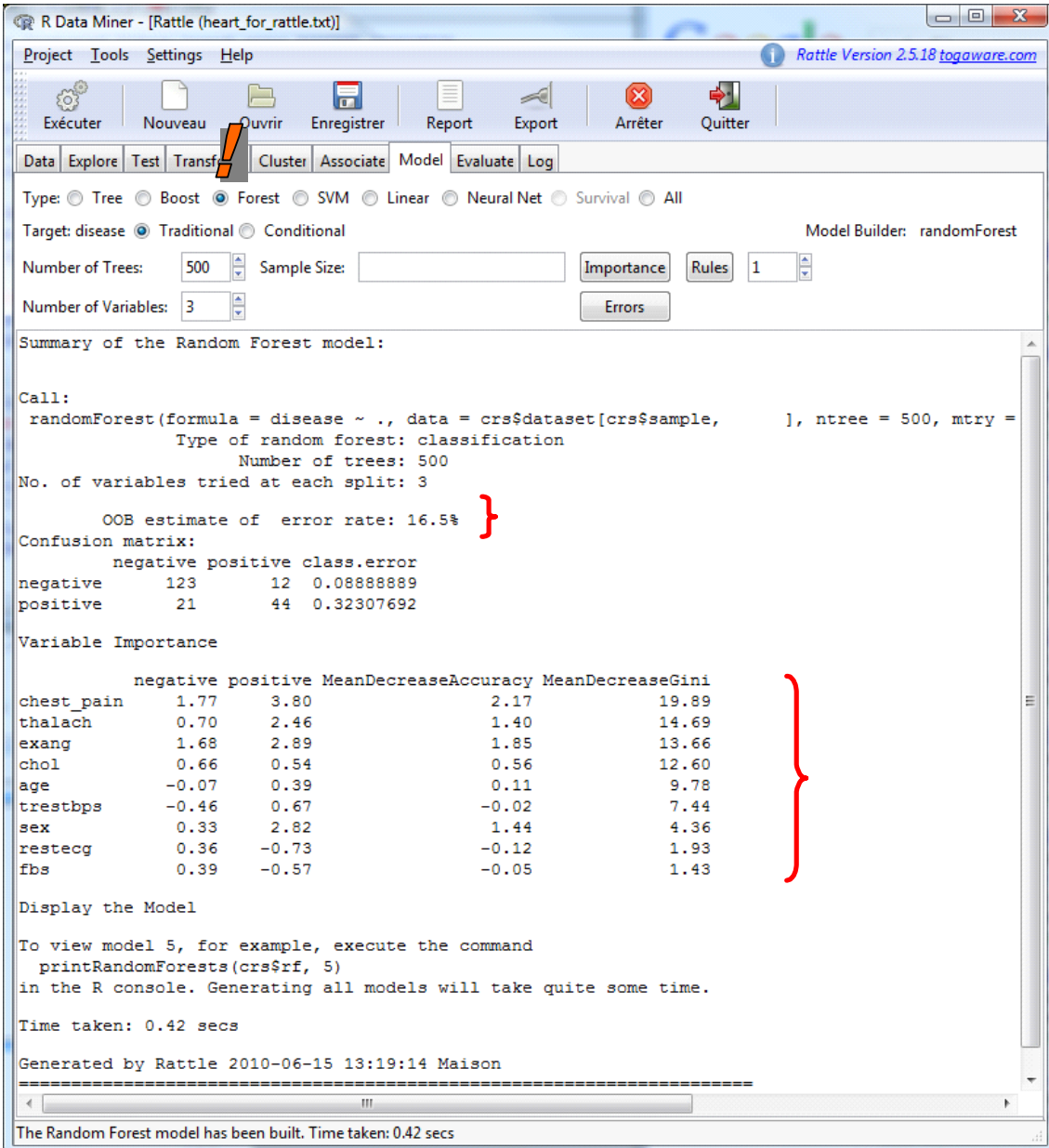
[1] 9 7 6 3 8 4 1 5 2

Generated by Rattle 2010-06-15 10:59:16 Maison
    
```

We can obtain also a graphical representation of the tree with the DRAW option.



About the **random forest** approach, rattle uses the `randomForest` command from the `randomForest` package. We obtain the following results with the default settings.



The screenshot shows the Rattle software interface. The main panel displays the following information:

Type: Tree Boost Forest SVM Linear Neural Net Survival All

Target: disease Traditional Conditional

Model Builder: randomForest

Number of Trees: 500 Sample Size: [] Importance Rules 1

Number of Variables: 3 Errors

Summary of the Random Forest model:

```
Call:
randomForest(formula = disease ~ ., data = crs$dataset[crs$sample, ], ntree = 500, mtry = 
              Type of random forest: classification
              Number of trees: 500
No. of variables tried at each split: 3
              OOB estimate of error rate: 16.5%
Confusion matrix:
      negative positive class.error
negative  123      12 0.08888889
positive   21      44 0.32307692
Variable Importance
      negative positive MeanDecreaseAccuracy MeanDecreaseGini
chest_pain  1.77    3.80           2.17           19.89
thalach    0.70    2.46           1.40           14.69
exang     1.68    2.89           1.85           13.66
chol      0.66    0.54           0.56           12.60
age      -0.07    0.39           0.11            9.78
trestbps -0.46    0.67          -0.02            7.44
sex       0.33    2.82           1.44            4.36
restecg   0.36   -0.73          -0.12            1.93
fbs       0.39   -0.57          -0.05            1.43
```

Display the Model

To view model 5, for example, execute the command

```
printRandomForests(crs$rf, 5)
```

in the R console. Generating all models will take quite some time.

Time taken: 0.42 secs

Generated by Rattle 2010-06-15 13:19:14 Maison

The Random Forest model has been built. Time taken: 0.42 secs

The OOB (out-of-bag) error estimation is 16.5%. We will compare this value to the one obtained on the test set below.

About the **logistic regression**, we use the `glm()` command. It automatically transforms the discrete predictors using dummy variables. We obtain the following results.

R Data Miner - [Rattle (heart_for_rattle.txt)]

Project Tools Settings Help Rattle Version 2.5.18 togaware.com

Exécuter Nouveau Ouvrir Enregistrer Report Export Arrêter Quitter

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: Tree Boost Forest SVM Linear Neural Net Survival All

Numeric Generalized Poisson Logistic Probit Multinomial Model Builder: glm (logit)

Plot

Summary of the Logistic Regression model (built using glm):

Call:
`glm(formula = disease ~ ., family = binomial(link = "logit"), data = crs$dataset[crs$sample,])`

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.5490	-0.4750	-0.2661	0.4889	2.8434

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	0.294096	4.036321	0.073	0.941916
age	-0.007997	0.033876	-0.236	0.813390
sexmale	1.189818	0.587423	2.025	0.042818 *
chest_painatyp_angina	-2.186895	0.575810	-3.798	0.000146 ***
chest_painnon_anginal	-1.586384	0.627181	-2.529	0.011426 *
chest_paintyp_angina	0.517539	0.923208	0.561	0.575079
trestbps	-0.004770	0.012892	-0.370	0.711408
chol	0.005226	0.003447	1.516	0.129502
fbst	1.251407	0.833696	1.501	0.133346
restecgnormal	1.098985	2.277026	0.483	0.629351
restecgst_t_wave_abnormality	0.426301	2.314418	0.184	0.853862
thalach	-0.023186	0.011624	-1.995	0.046071 *
exangyes	1.874358	0.494388	3.791	0.000150 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 252.23 on 199 degrees of freedom
 Residual deviance: 144.08 on 187 degrees of freedom
 AIC: 170.08

Number of Fisher Scoring iterations: 6

Log likelihood: -72.039 (13 df)
 Null/Residual deviance difference: 108.153 (12 df)
 Chi-square p-value: 0.00000000
 Pseudo R-Square (optimistic): 0.70482172

3.6 Measuring the generalization performance

Last step of our analysis, we want to evaluate the performances of the classifiers on the test sample (30% of the whole dataset).

We activate the "Evaluate" tab. First, we want to obtain the confusion matrix and the associated error rate. We select the "Error Matrix" option. For the "Data" item, we must select the "Testing" option. Only the models learned into the "Model" tab are available here.

We click on the EXECUTER menu. We observe that the logistic regression is the better here with a test error rate equal to 18.18%.

We note also that the OOB error rate (16.5%) seems underestimate the error rate for the random forest (20.45% on the test set). But, because the test set size is small, and the test error rate being also an estimation of the "true" error rate, we consider with many cautions this result.

R Data Miner - [Rattle (heart_for_rattle.txt)]

Project Tools Settings Help

Rattle Version 2.5.18 togaware.com

Executer Nouveau Ouvrir Enregistrer Report Export Arrêter Quitter

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: Error Matrix Risk Cost Curve Hand Lift ROC Precision Sensitivity Prv Ob Score

Model: Tree Boost Forest SVM Linear Neural Net Survival KMeans HClust

Data: Training Validation Testing CSV File (Aucun) R Dataset

Risk Variable: Report: Class Probability Include: Identifiers All

Error matrix for the Decision Tree model on heart_for_rattle.txt [test] (counts):

	Actual	
Predicted	negative	positive
negative	23	6
positive	3	12

Error matrix for the Decision Tree model on heart_for_rattle.txt [test] (%):

	Actual	
Predicted	negative	positive
negative	52	14
positive	7	27

Overall error: 0.2045455

Generated by Rattle 2010-06-15 13:48:39 Maison

Error matrix for the Random Forest model on heart_for_rattle.txt [test] (counts):

	Actual	
Predicted	negative	positive
negative	25	8
positive	1	10

Error matrix for the Random Forest model on heart_for_rattle.txt [test] (%):

	Actual	
Predicted	negative	positive
negative	57	18
positive	2	23

Overall error: 0.2045455

Generated by Rattle 2010-06-15 13:48:39 Maison

Error matrix for the Linear model on heart_for_rattle.txt [test] (counts):

	Actual	
Predicted	negative	positive
negative	24	6
positive	2	12

Error matrix for the Linear model on heart_for_rattle.txt [test] (%):

	Actual	
Predicted	negative	positive
negative	55	14
positive	5	27

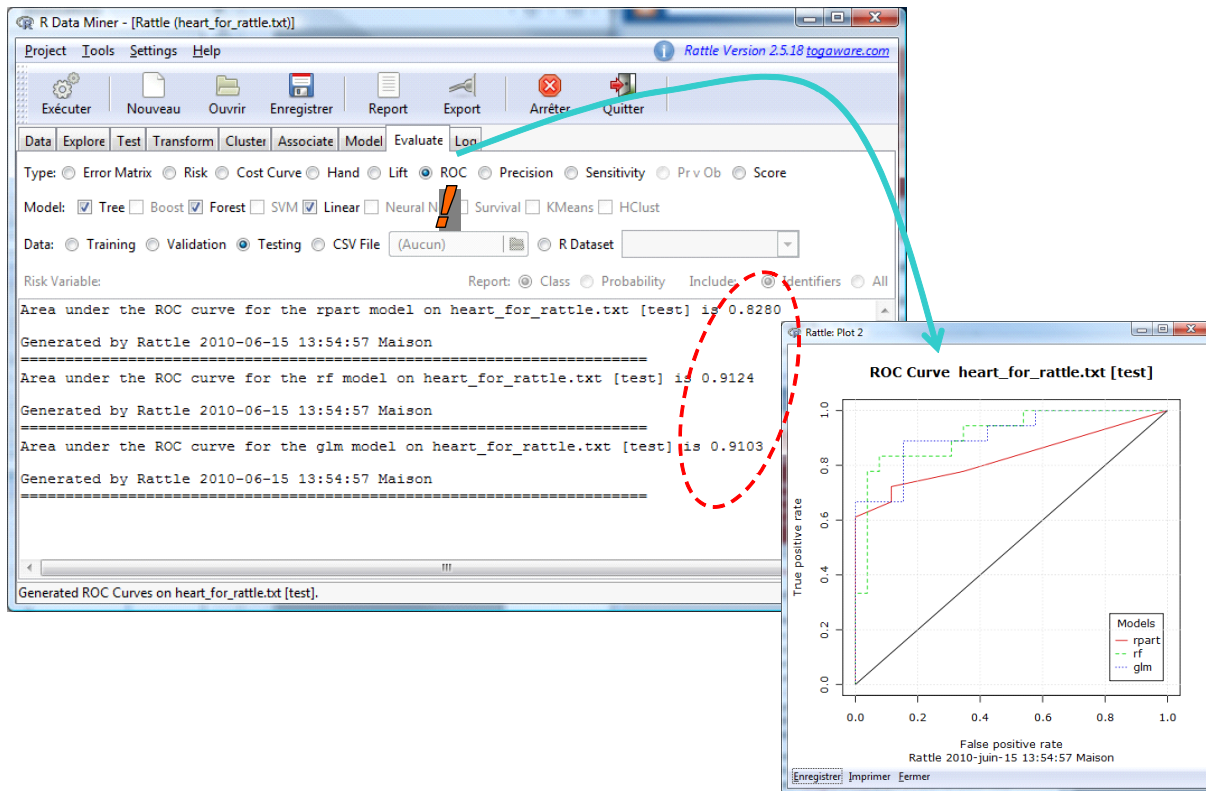
Overall error: 0.1818182

Generated by Rattle 2010-06-15 13:48:39 Maison

Generated Error Matrix.

Actually, the error rate is not a good criterion here. We note that the differences between the methods are based only on one misclassified instance. In our context, it is perhaps more interesting to use the ROC curve which highlights the ability of the methods to assign higher score to the positive instances compared with the negative ones (see <http://data-mining-tutorials.blogspot.com/2008/11/roc-curve-for-classifier-comparison.html> or <http://data-mining-tutorials.blogspot.com/2008/10/computing-roc-curve.html>).

We select the ROC option under rattle.



According to the AUC criterion, the decision tree is definitely the worst compared with the two other classifiers, which are similar in terms of performance. It is not surprising. We know that the decision tree is not well adapted to the [scoring](#) process.

3.7 R commands associated to the treatments

```

R Data Miner - [Rattle (heart_for_rattle.txt)]
Project Tools Settings Help
Rattle Version 2.5.18 togaware.com
Exécuter Nouveau Ouvrir Enregistrer Report Export Arrêter Quitter
Data Explore Test Transform Cluster Associate Model Evaluate Log
[ ] Export Comments [ ] Rename Rattle Variables: From crs$ to MY
# little effort the log can be used to score a new dataset. The logical variable
# 'building' is used to toggle between generating transformations, as when building
# a model, and simply using the transformations, as when scoring a dataset.

building <- TRUE
scoring <- ! building

# The colorspace package is used to generate the colours used in plots, if available.

library(colorspace)

=====
# Rattle timestamp: 2010-06-15 14:34:33 i386-pc-mingw32

# Load the data.

crs$dataset <- read.csv("file:///D:/DataMining/Databases_for_mining/benchmark_datasets/Heart/heart_for_rattle.txt", sep="\t", na.strings=c(".", "
=====
# Rattle timestamp: 2010-06-15 14:34:34 i386-pc-mingw32

# Note the ...
    
```

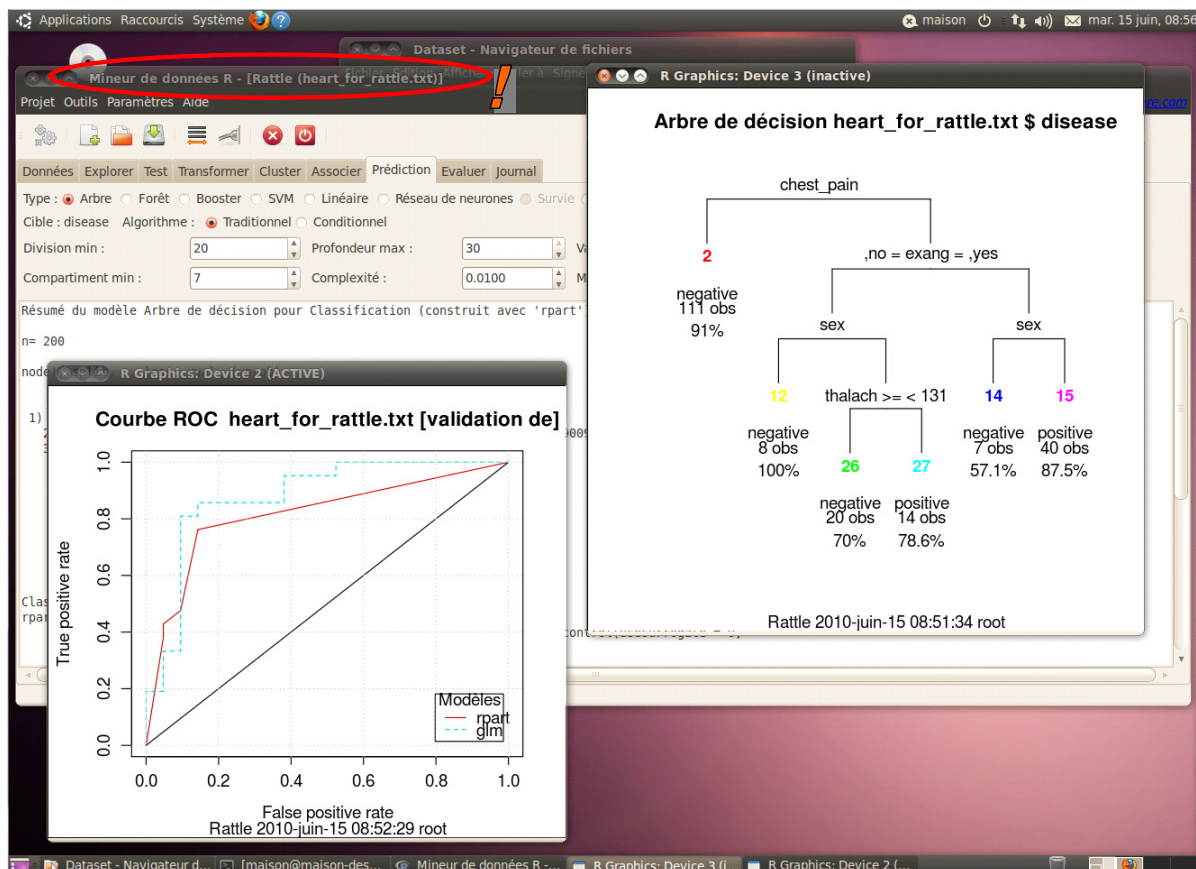
One of the main criticisms which we make for the software driven by menu is that once the process is finalized, when we close the software, we have no recollection of the sequence of operations we performed. In the next working session, it is complicated to reproduce them as before. It is necessary to have an excellent memory, or to have taken care of noting all that we made.

Rattle allows to overtake this drawback by translating all the operations (corresponding to a click on the EXECUTER menu) performed by the user in a sequence of R commands. We can visualize them in the "Log" tab. We can store these commands (and the comments) into a file. In the next working session, it is very easy to perform the same data processing by loading these commands.

4 Rattle under Linux (Ubuntu)

The installation of the Rattle package under Linux is not easy. We must follow carefully the description available on the website. In case of problem, a troubleshooting procedure is proposed. This is the one that I used (see http://datamining.togaware.com/survivor/Install_GNU_Linux.html).

When the installation is finalized, Rattle works properly under Linux (Ubuntu) as we see below.



5 Conclusion

In this tutorial, we showed that it was possible to use R without knowledge about its programming language with the help of the rattle package. This package is rather specialized about the data mining methods. For the statisticians, there are other packages such as "[R Commander](#)".