# 1  Theme

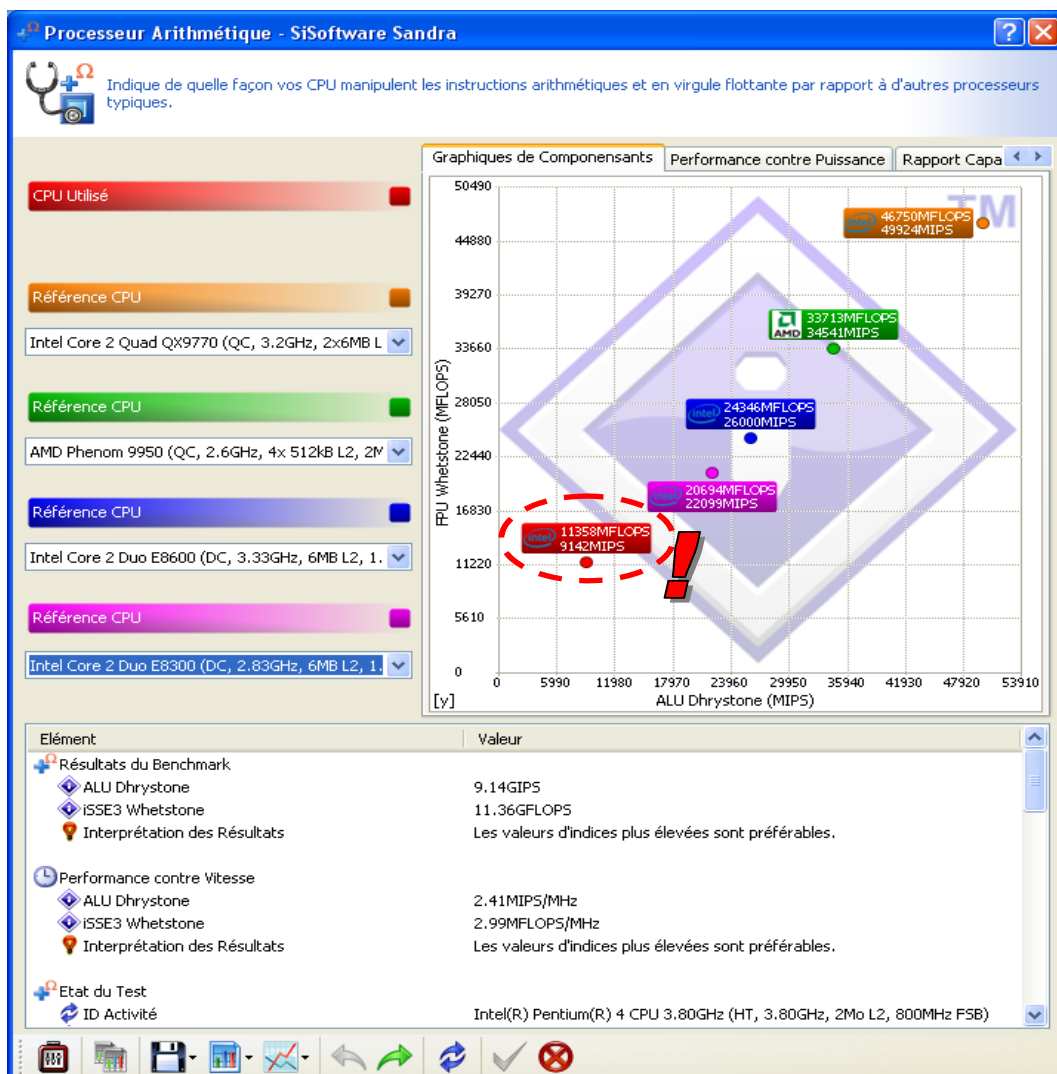**Comparing the quickness and memory occupation of various free implementation of the C4.5 (or assimilated) decision tree algorithm on large dataset.**

**2011/10/29 - The results are updated from a new computer (Q9400 core 2 quad processor) and a new operating system (Windows 7 - 64 bit) – See section 5.**

Dealing with large dataset is one of the most important challenges of the Data Mining. In this context, it is interesting to analyze and to compare the performances of various free implementations of the learning methods, especially the computation time and the memory occupation. Most of the programs download all the dataset into memory. The main bottleneck is the available memory.

In this tutorial, we compare the performance of several implementations of the C4.5 algorithm (Quinlan, 1993) when processing a file containing 500,000 observations and 22 variables.

Our main criteria are memory occupation and computation time. So that everyone can replicate the experience and compare the results, here are the characteristics of our machine: a Pentium 4 3.8 GHz with 2 GB of RAM. We have measured the performances of our computer using the Lite version of the SISOFTWARE SANDRA (http://www.sisoftware.net/). We obtain:

We compare the performances of the following programs in this tutorial.

| Software | Version | URL |
|----------|---------|-----|
| KNIME | 1.3.5 | http://www.knime.org/index.html |
| ORANGE | 1.0b2 | http://www.ailab.si/orange/ |
| R (package rpart) | 2.6.0 | http://www.r-project.org/ |
| RAPIDMINER (YALE) | Community Edition | http://rapid-i.com/ |
| SIPINA | Research | http://eric.univ-lyon2.fr/~ricco/sipina.html |
| TANAGRA | 1.4.27 | http://eric.univ-lyon2.fr/~ricco/tanagra/ |
| WEKA | 3.5.6 | http://www.cs.waikato.ac.nz/ml/weka/ |

# 2   The WAVE dataset

Our data file is **WAVE500K .ZIP** (http://eric.univ-lyon2.fr/~ricco/tanagra/fichiers/wave500k.zip). It is the well-known artificial dataset described in the CART book (Breiman et al., 1984). We have generated a dataset with 500.000 observations. The class attribute has 3 values, there are 21 continuous predictors.

We deal primarily with ARFF WEKA file format. The majority of data mining software can handle this format. Otherwise, text file format with tabulation separator will be used.
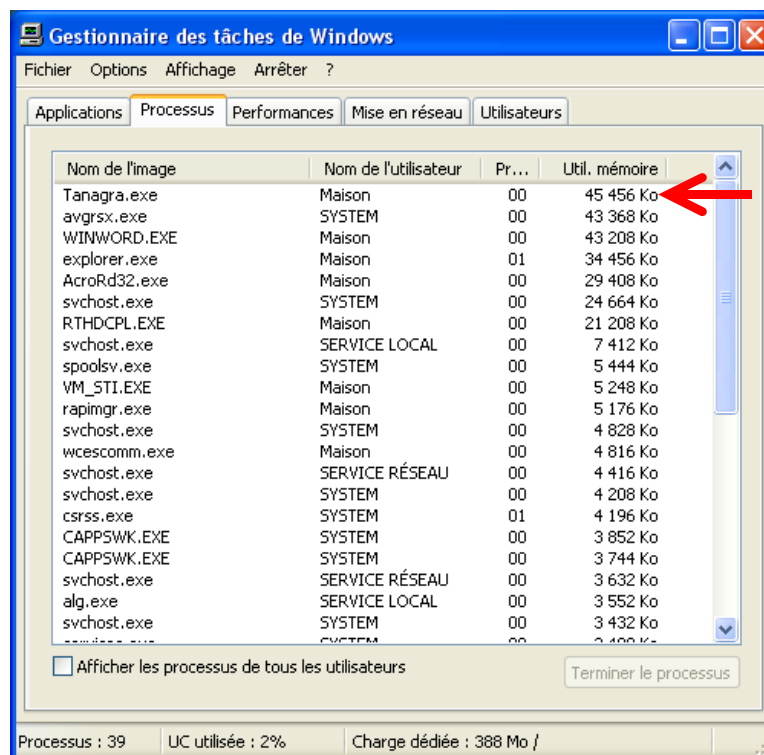
# 3   Comparison of performances



**Figure 1 – Measuring the memory occupation under Windows OS**
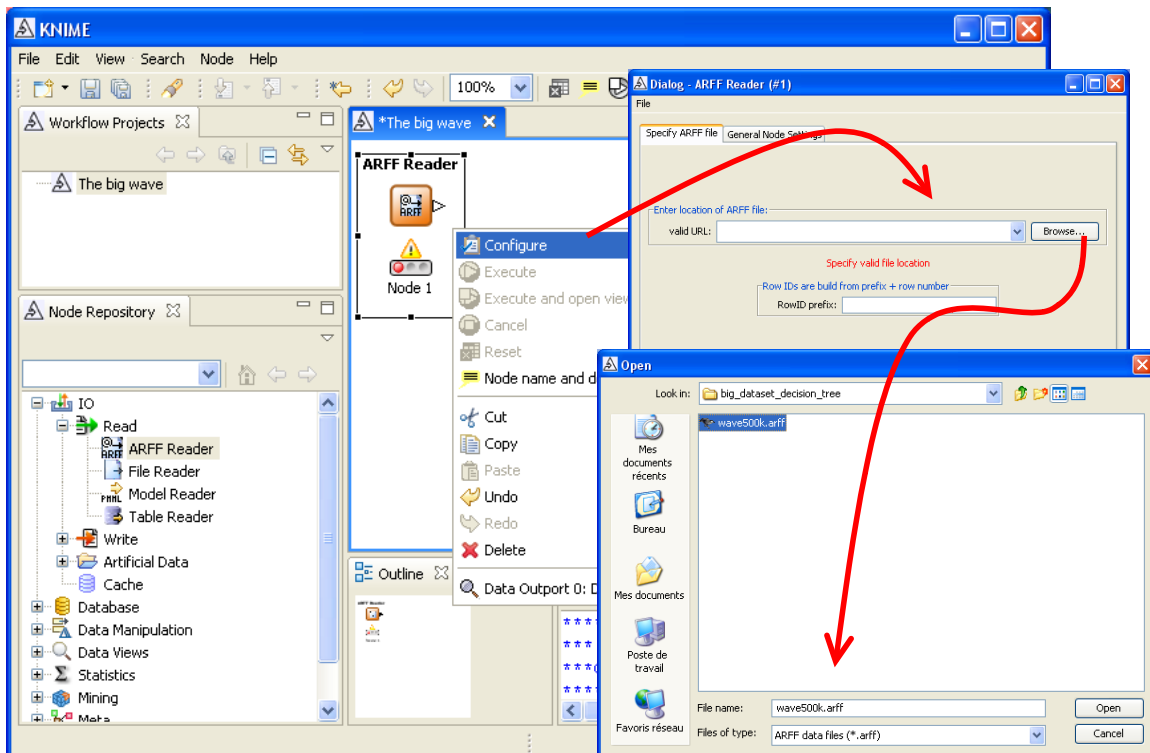
We utilize the indications of Windows OS about the memory occupation (Figure 1). We measure the memory occupation when the program is launched, when the dataset is loaded, during the learning phase, after the learning phase. The maximum occupation during the learning phase is the most important indicator. This is the main bottleneck of the whole process.

About the computation time, some programs provide it. For the other cases, we use a chronograph. Indeed, we try to obtain only an order of magnitude of the computation time. It is not necessary to measure a very precise value.

## 3.1   KNIME

KNIME (Konstanz Information Miner -- http://www.knime.org/) describes the treatments as a succession of operations, represented by a workflow.
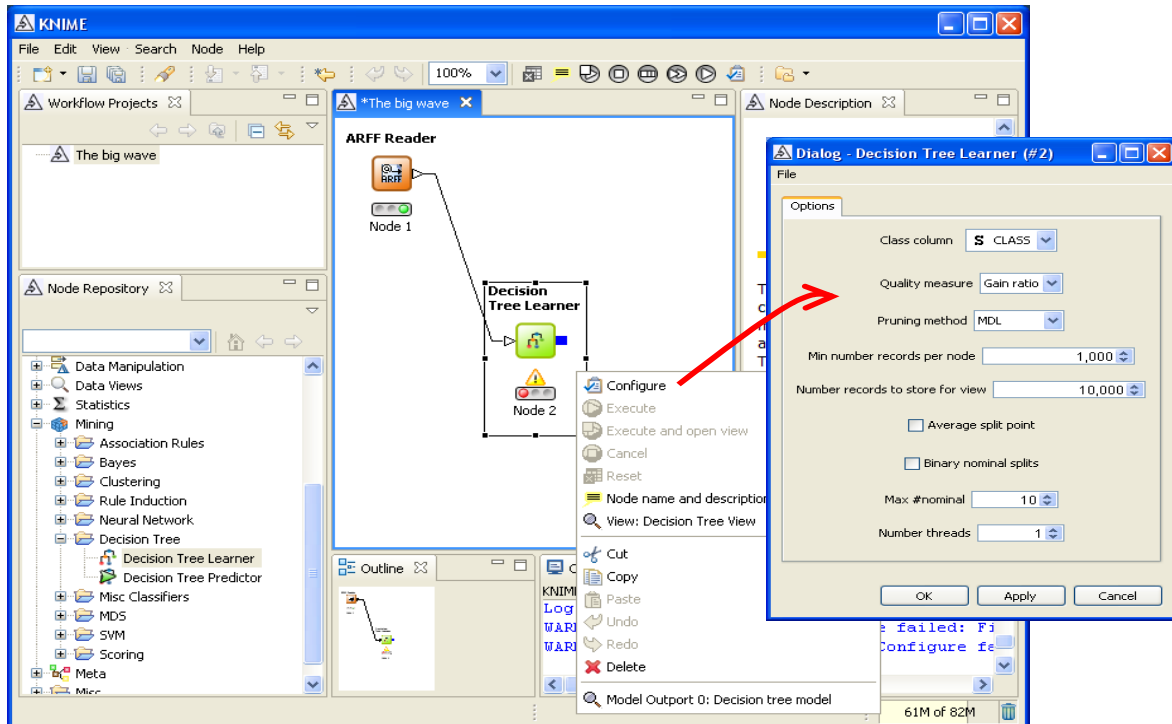
**KNIME and the JRE**. When we launch the software, the JRE (Java Runtime Environment) is started and the global memory occupation is 92.6 MB.
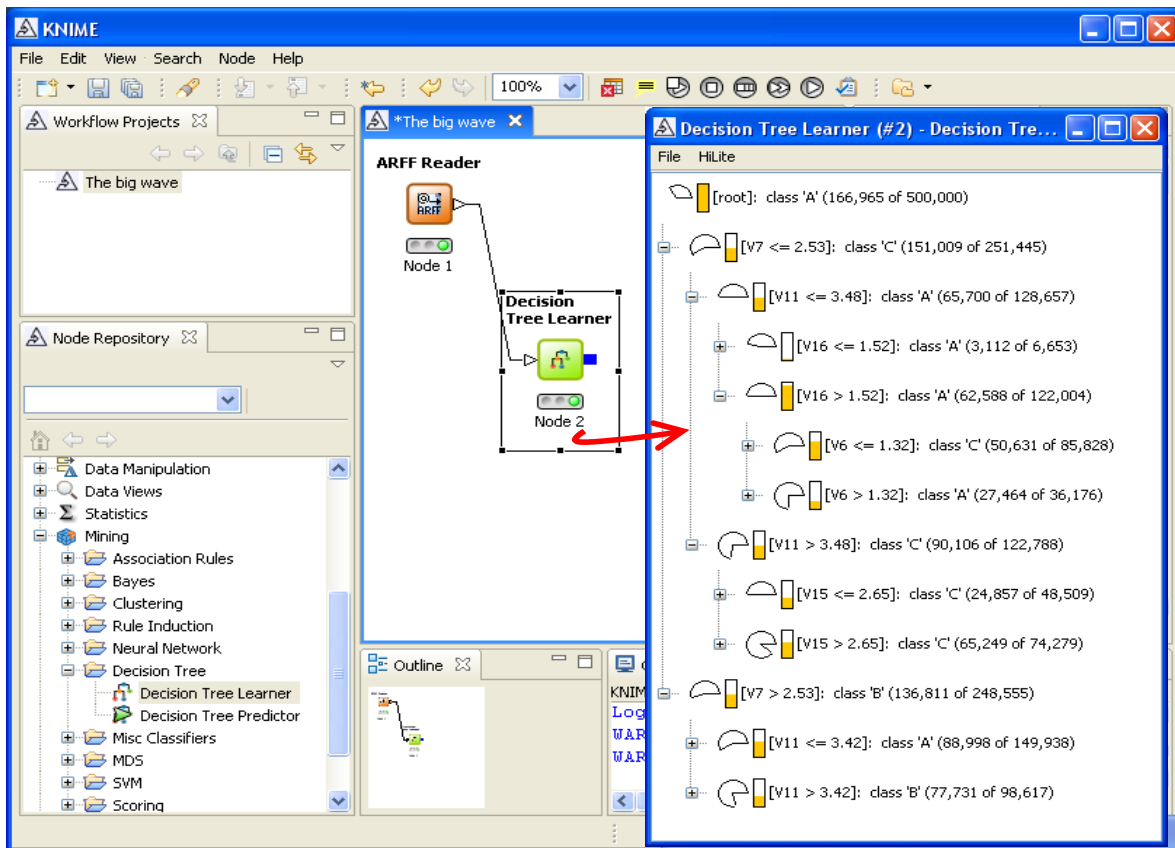


We create a new project (FILE / NEW). We add the ARFF READER component into the workflow and we click on the CONFIGURE menu in order to select the data file. After the loading, the memory occupation becomes 160.4 MB.

We insert the decision tree LEARNER in the workflow, we see on the panel at the right a description of the method and its parameters. We connect the component ARFF READER, then we click on the menu CONFIGURE in order to define the parameters.

It is not really the C4.5 method, but it is quite similar. We set the minimum number of examples per node at 1000. We set the number of thread used to 1 in order to make the results comparable. We note that Knime is the only one software that can use several threads during the induction of the decision tree.

We launch the calculation by clicking on the EXECUTE AND OPEN VIEW menu. Knime does not give indications about the tree structure (number of nodes and leaves). The processing time is 270 seconds; the maximum memory occupation is 245.8 MB during the learning phase.
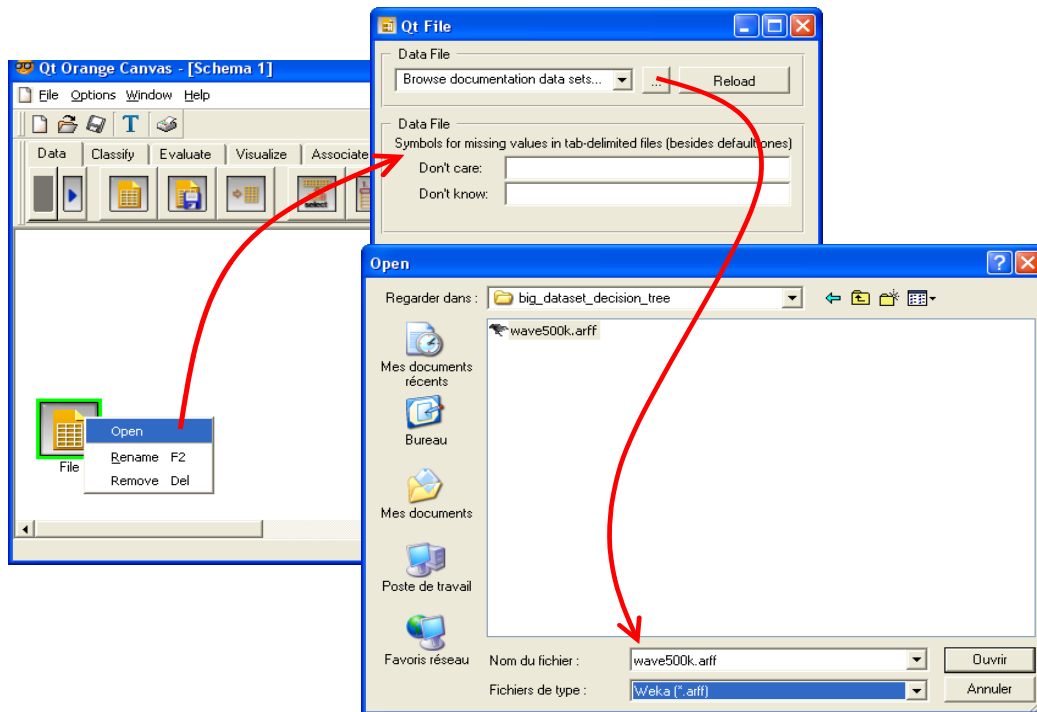


It is possible to save the workflow. But, it seems that the software tries to store all the available information, including the dataset. The processing time can be very important for our analysis.
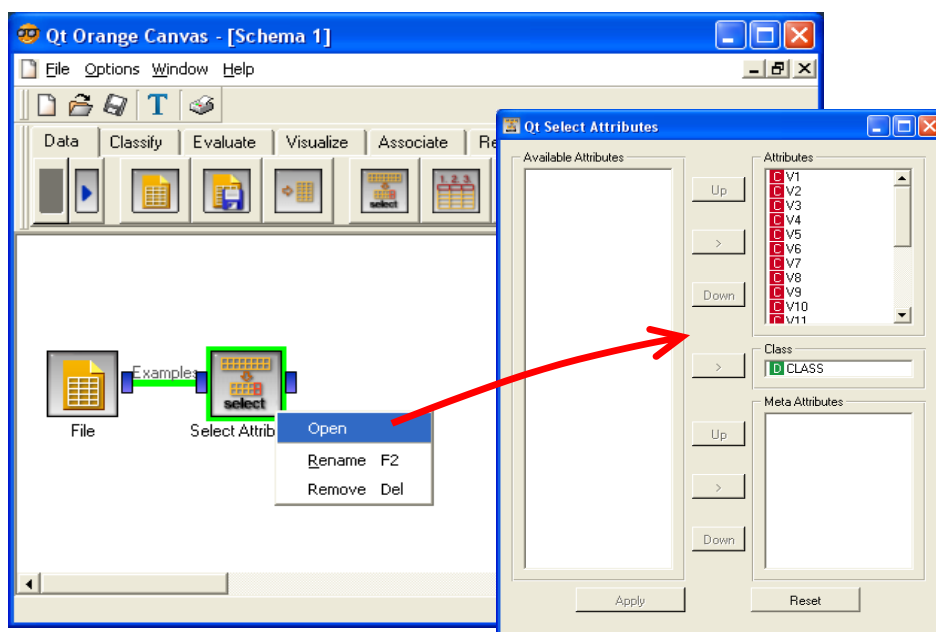
## 3.2   ORANGE

ORANGE can be used through standard programming language (python) or visual programming. We use this last framework in this tutorial.
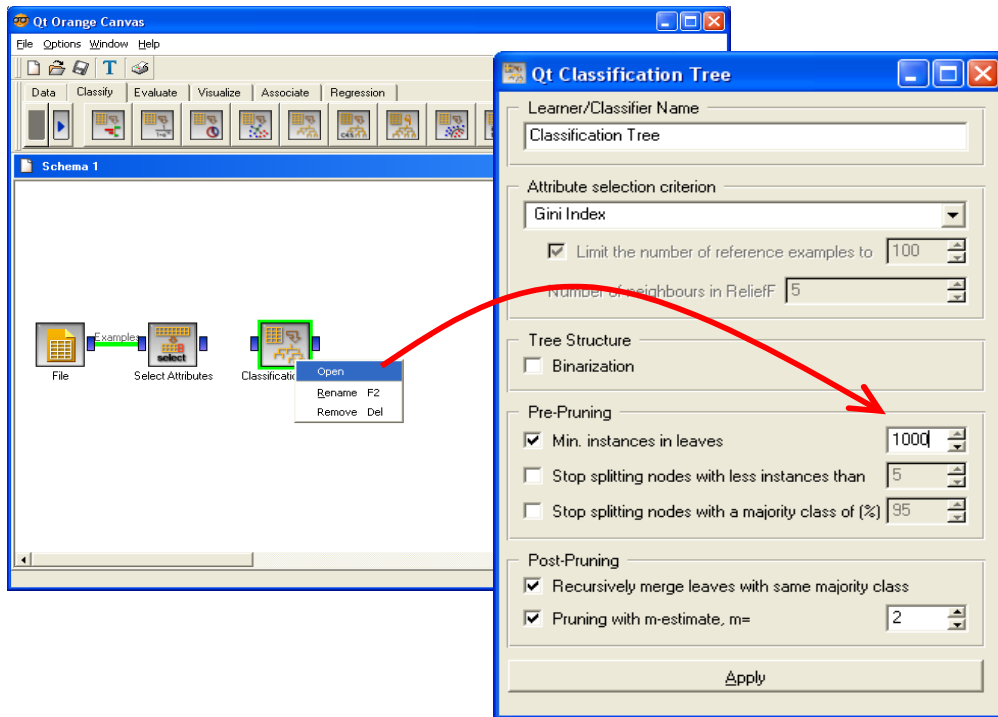
When we launch ORANGE, the **Python Runtime Environment** is started. The global memory occupation is 24.9 MB. We have a familiar interface. The treatments are defined in a schema. The components are available into the tool bar. We add the FILE (DATA tab) into the schema. We click on the OPEN menu and we select the data file. The loading is automatically started; there is no specific menu for launching the operation.
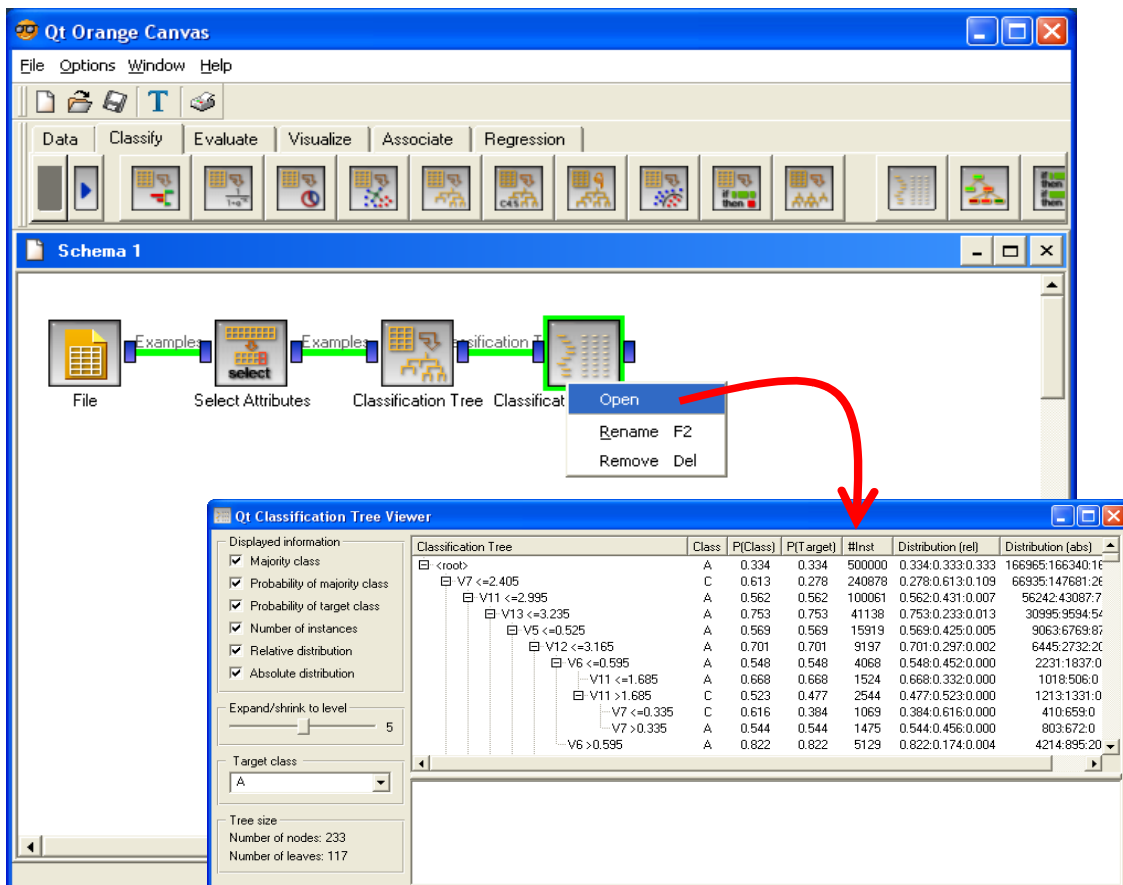
The computation time is 90 seconds. The memory occupation becomes 259.5 MB. We add the SELECT ATTRIBUTE component in order to define the position of the variables.

We can insert now the CLASSIFICATION TREE (CLASSIFY tab) component. We configure it **before** we make the connection. We click on the OPEN MENU. We set the parameters of the algorithm.



We insert the CLASSIFICATION TREE VIEWER (CLASSIFY tab) into the schema. We make the connection between this viewer and the learner. Finally, we make the connection between SELECT ATTRIBUTES and CLASSIFICATION TREE. The computation is automatically launched.

The processing time is 130 seconds; the memory occupation is 795.7 MB. The tree includes 117 leaves.

## 3.3   The R Software with the RPART package

R works with scripts. We use the **system.time(.)** command to obtain the processing time.

```
R RGui - [R Console]
R  Fichier  Edition  Voir  Misc  Packages  Fenêtres  Aide

R version 2.6.0 (2007-10-03)
Copyright (C) 2007 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.

> #chargement des données
> setwd("D:/DataMining/Databases_for_mining/comparison_TOW/big_dataset_decision_tree")
> system.time(donnees <- read.table(file="wave500k.txt",dec=".",header=T))
utilisateur     système      écoulé
     24.17        0.17        24.66
> #construction de l'arbre
> library(rpart)
> parametres <- rpart.control(xval=0,minsplit=1000,minbucket=1000,maxcompete=0,maxsurrogate=0,cp=0)
> system.time(modele <- rpart(CLASS ~., data = donnees, method = "class", control = parametres))
utilisateur     système      écoulé
    154.62        1.67       157.09
```

At the beginning, the memory occupation is 18.8 MB. We use the text file format with R, we load the WAVE500K.TXT. The computation is 24 seconds; the memory occupation becomes 184.1MB.
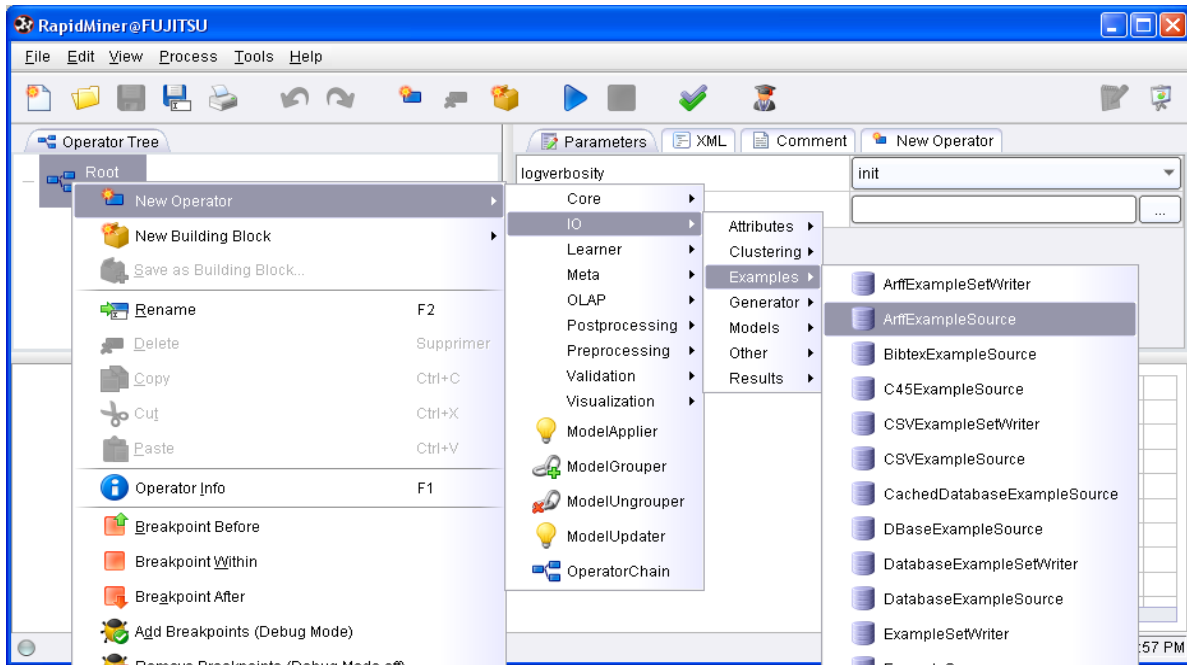
The RPART procedure is rather similar to CART method. But we set the parameters in order to produce a tree similar to the others programs i.e. MINSPLIT = 1000 observations; MINBUCKET = 1000 observations; CP = 0. There is not post-pruning here. But, it does not matter. We know that this part of the processing is fast in the C4.5 approach.

The processing time for the tree construction is 157 seconds. The memory occupation becomes 718.9 MB.
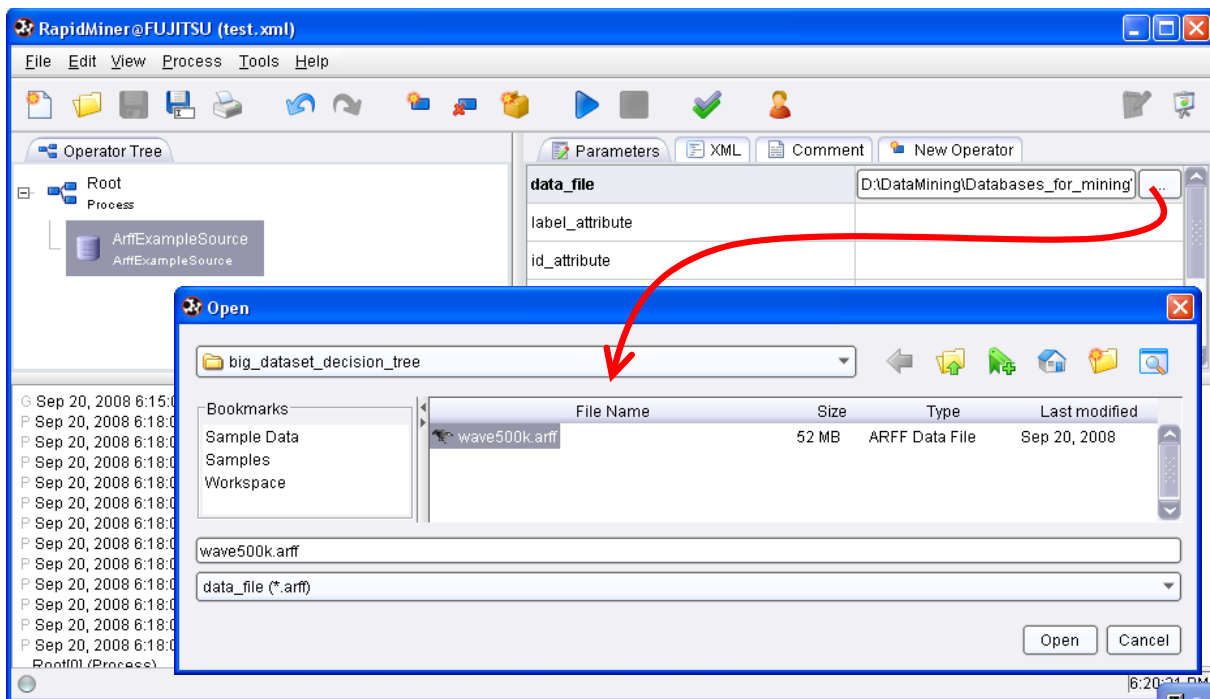
## 3.4   RAPIDMINER (formerly YALE)

When we launch RAPIDMINER, we create a new project by clicking on FILE / NEW. An empty OPERATOR TREE is available.

We add an ARFFEXAMPLESOURCE component into the tree, in order to load the dataset.
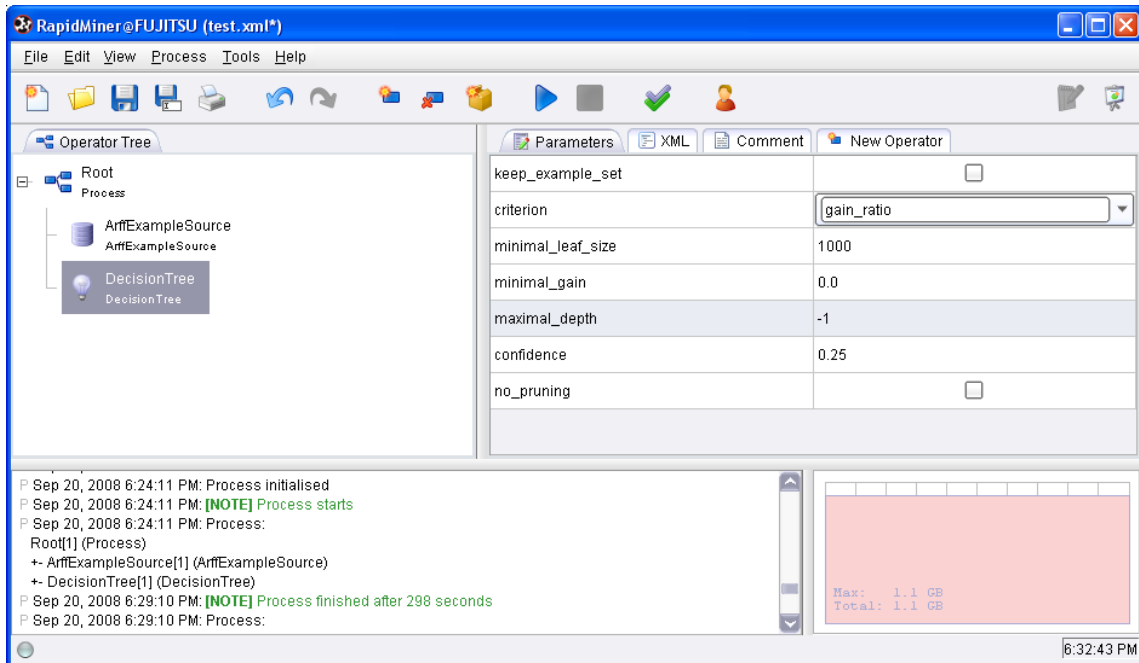
The component appears into the operator tree, we set the right parameters by selecting our data file. We click on the PLAY button of the tool bar. The computation time for the data importation is 7 seconds. The allocated memory increases from 136.6MB to 228.1MB.



We must now insert the "decision tree" into the diagram. We operate from the menu NEW OPERATOR on the root of the tree. We use the DECISIONTREE operator. We use parameters similar to the other programs.
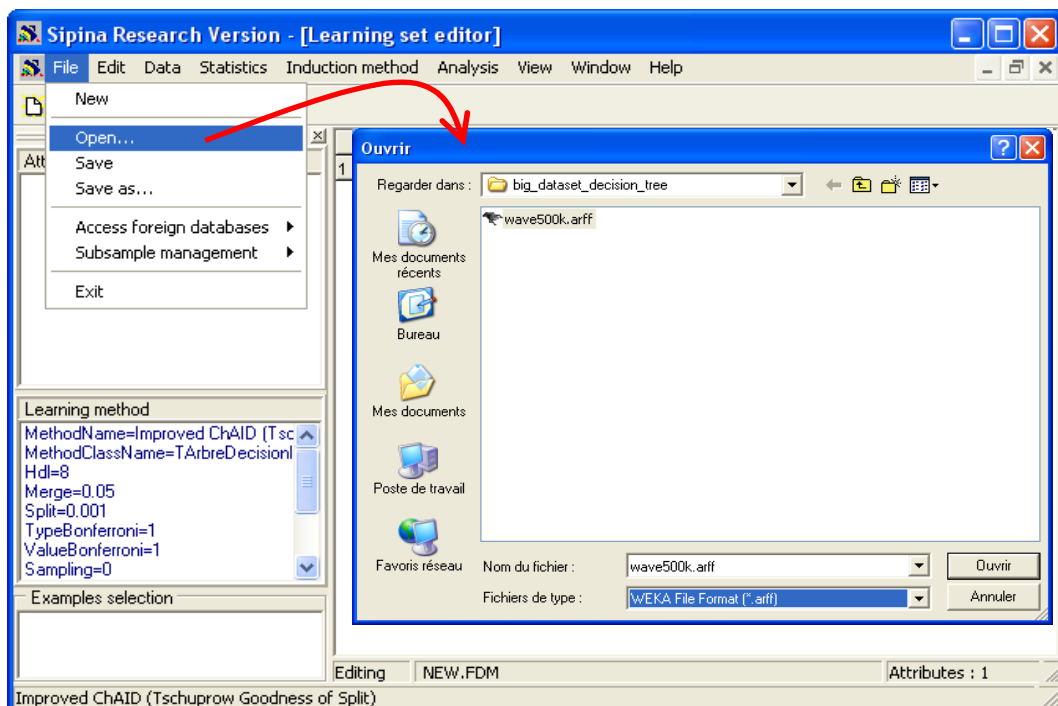
The PLAY button starts again all the computations. We obtain the decision tree after 298 seconds.

Curiously, RAPIDMINER needs much memory: 1274.7MB were necessary during the calculations. It seems we must take with caution this value.
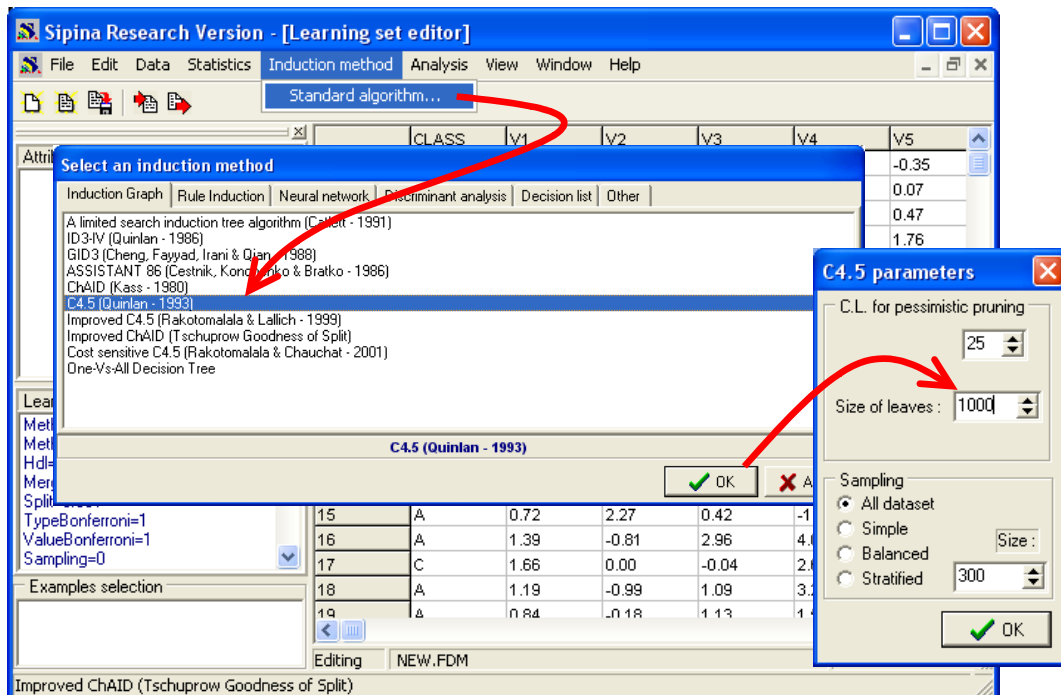
## 3.5   SIPINA

SIPINA is one of my old projects, specialized to decision tree induction. Unlike the other programs, SIPINA proposes a graphical display of the tree, and above all, possesses interactive functionalities.

This advantage is also a drawback because SIPINA must store much information on each node of the tree: the characteristics of split for the other variables, distributions, list of examples, etc. Thus, the memory occupation will be inevitably high compared to the other programs.
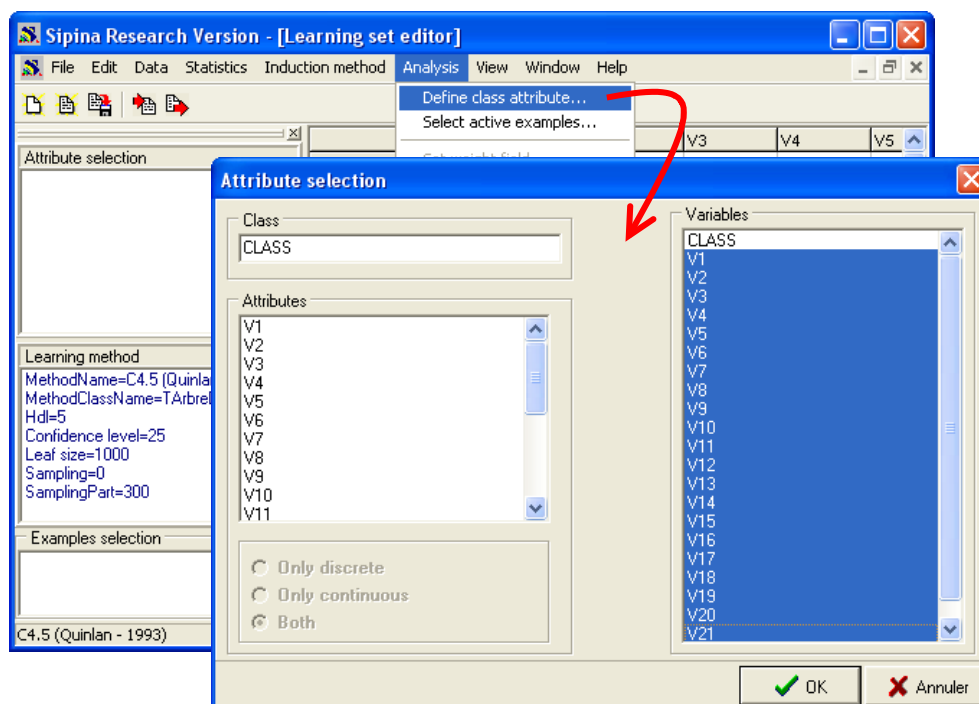
The memory occupation is 7.8 MB when we launch SIPINA. It becomes 67.1 MB after the dataset loading. The computation time is 25 seconds.
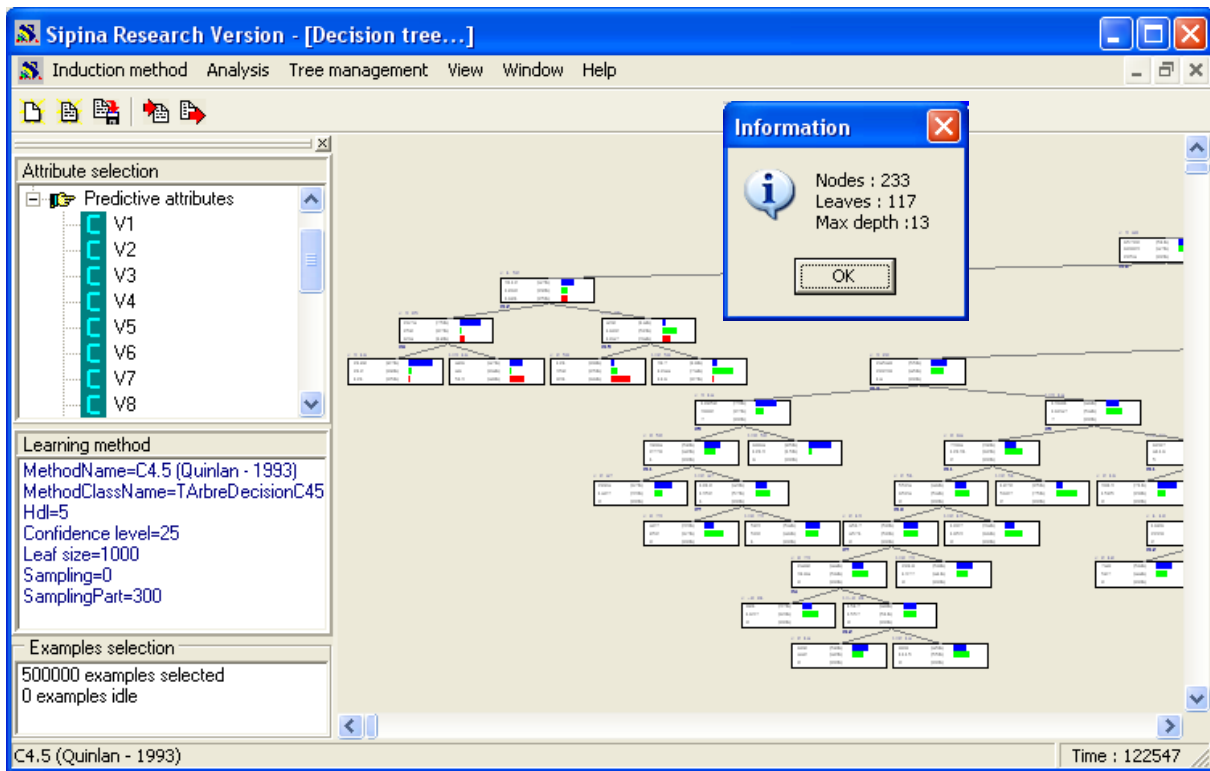
We select the C4.5 algorithm by clicking on the INDUCTION METHOD / STANDARD ALGORITHM menu. We set the minimal size of leaves to 1000.
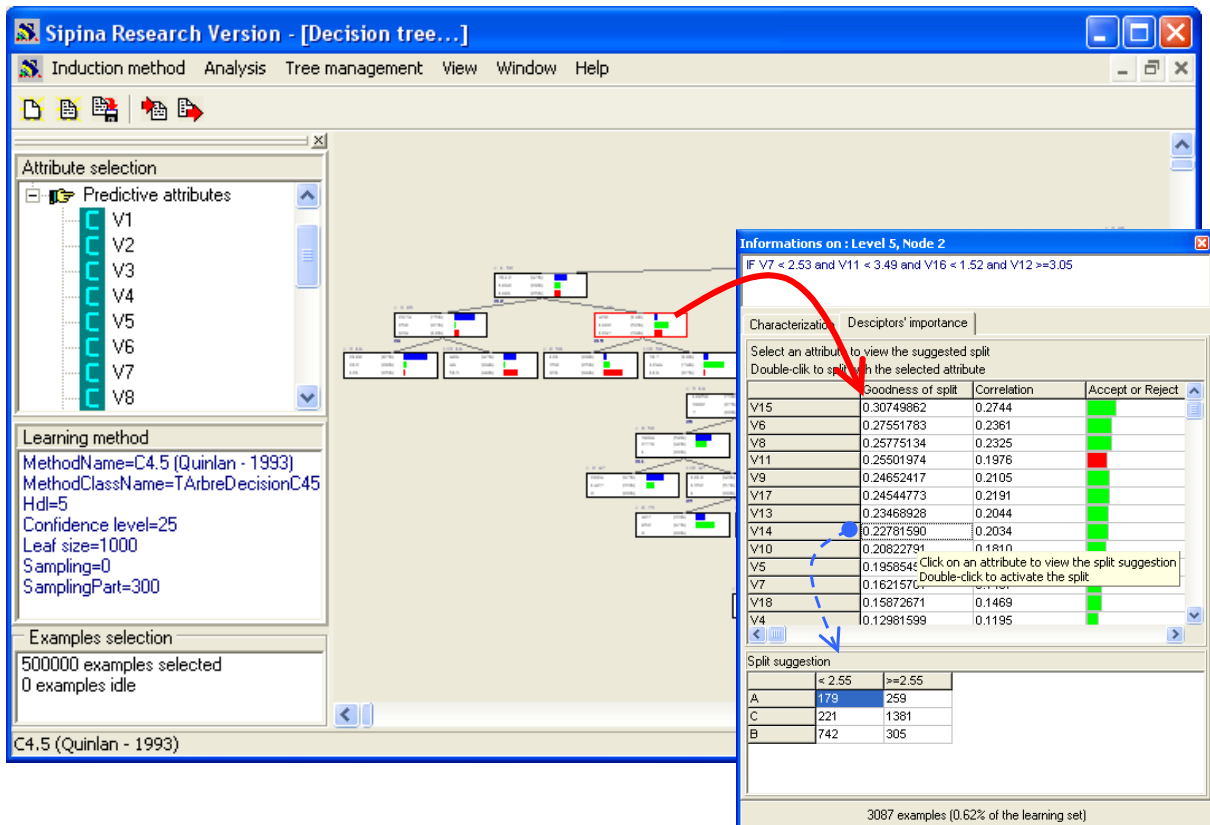


We must specify the role of the variables. We click on the ANALYSIS / DEFINE CLASS ATTRIBUTE menu. We set CLASS as TARGET, and the other variables (V1 … V21) as INPUT.



We can launch the learning phase. The calculation time is 122 seconds. The memory occupation becomes 539.9 MB for a tree with 117 leaves.
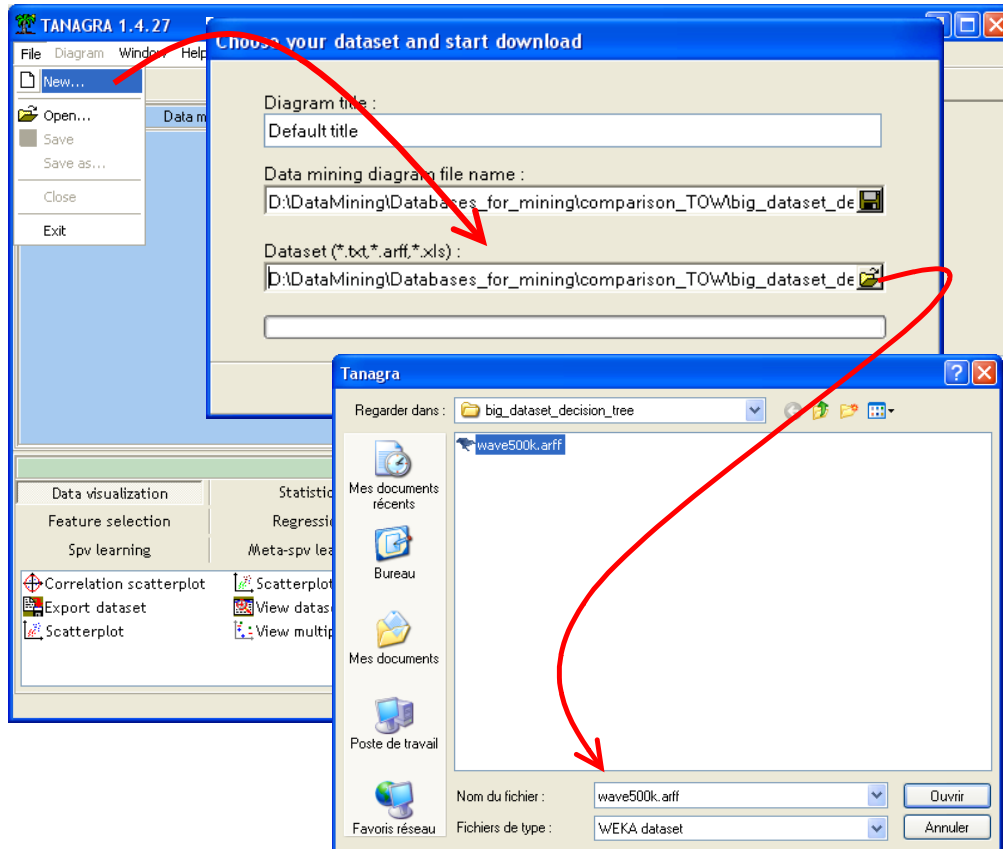
The counterpart of this large memory occupation is that we can explore deeply each node. In the screenshot below, we observe the splitting performance for each variable and, when we select a variable, we see the class distribution on the subsequent leaves.
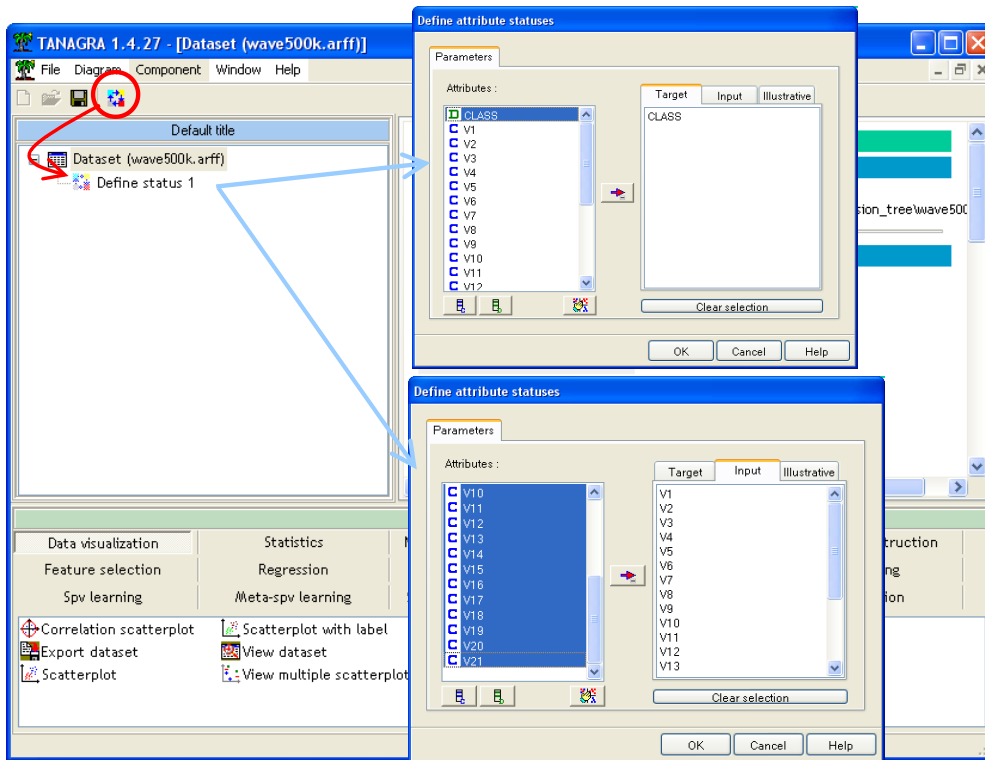
## 3.6   TANAGRA

TANAGRA is my current project. Unlike SIPINA, the implementation of the tree is widely modified. The interactive functionalities are anymore available, but the computation time and the memory occupation are optimized.

When we launch TANAGRA, the memory occupation is 7 MB. We click on the FILE / NEW menu in order to create a new diagram. We import the WAVE500K.ARFF data file.
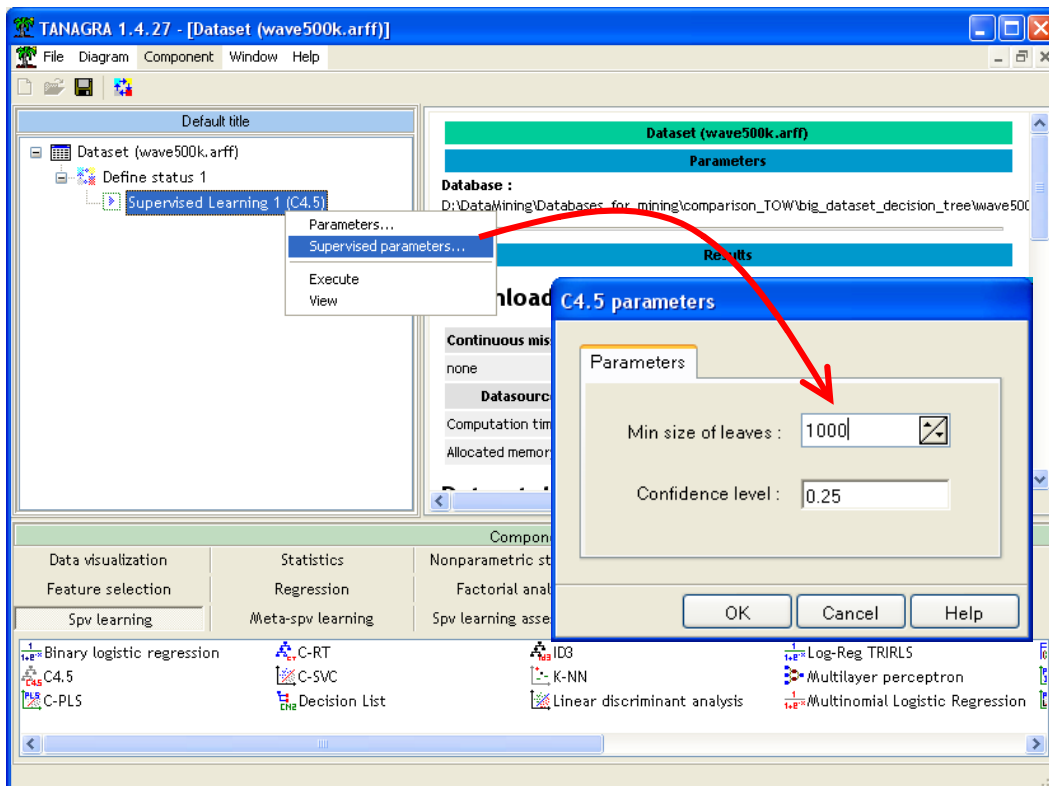


The processing time is 11 seconds; the memory occupation becomes 53.1 MB.

Before the learning phase, we must specify the role of each variable. We use the DEFINE STATUS component available into the tool bar. We set CLASS as TARGET, the other columns as INPUT (V1…V21).
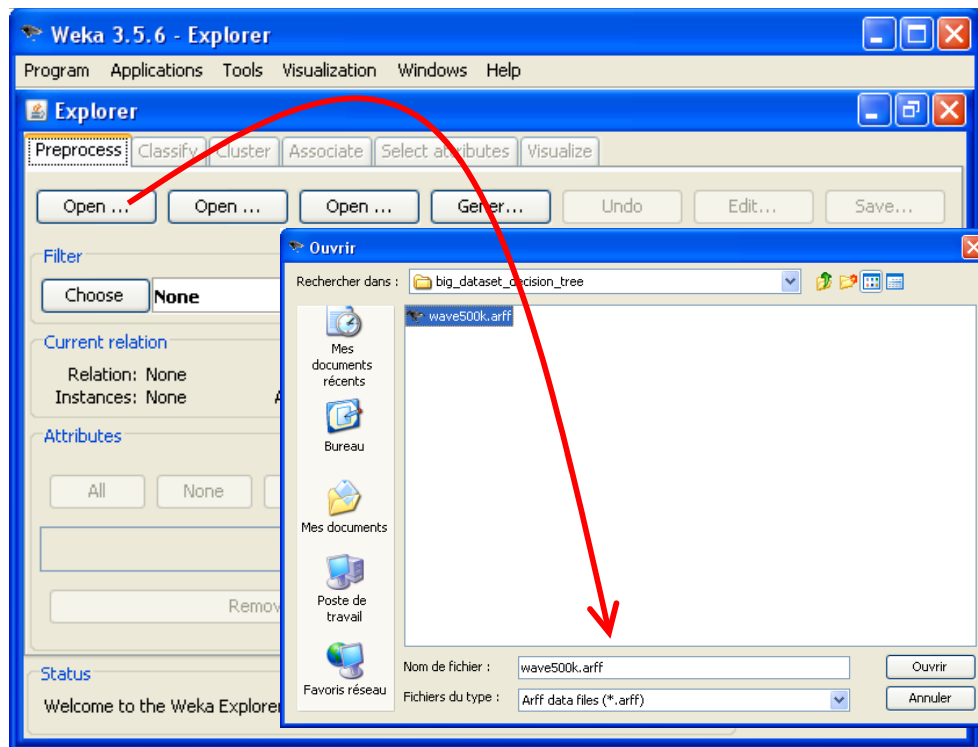
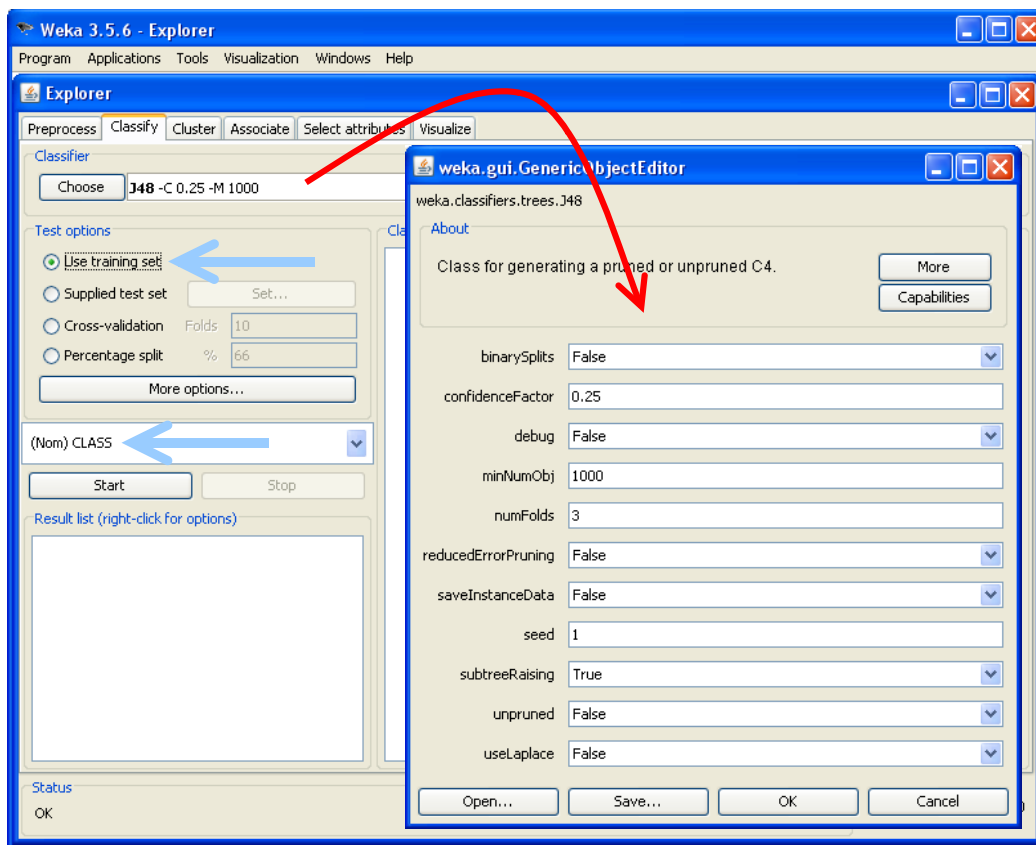We can now insert the C4.5 component and set the adequate parameters.



The contextual menu VIEW launches the calculations. The processing time is 33 seconds. The maximal memory occupation during the computation is 121.6 MB, it is 73.5MB when the tree is finished. The decision tree includes 117 leaves.

## 3.7   WEKA

We use the EXPLORER module for WEKA. We click on the OPEN button of the PREPROCESS tab. We select the WAVE500K.ARFF data file.



The processing time is 10 seconds; the memory occupation becomes 253.2 MB.

When the dataset is loaded, we select the CLASSIFY tab. We choose the J48 decision tree algorithm, we set the adequate parameters. We can now click on the START button.

The processing time is 338 seconds. The memory occupation becomes 699.6 MB. The tree includes 123 leaves.

### 3.8   Summary

The results are summarized into the following table:

| Program | Computation time (seconds) | | Memory occupation (MB) | | | |
|---|---|---|---|---|---|---|
| | Data Importation | Tree induction | After launch | After importation | Max during treatment | After induction |
| KNIME | 47 | 270 | 92.6 | 160.4 | 245.8 | 245.8 |
| ORANGE | 90 | 130 | 24.9 | 259.5 | 795.7 | 795.7 |
| R (package rpart) | 24 | 157 | 18.8 | 184.1 | 718.9 | 718.9 |
| RAPIDMINER | 7 | 298 | 136.3 | 228.1 | 1274.4 | 1274.4 |
| SIPINA | 25 | 122 | 7.8 | 67.1 | 539.9 | 539.9 |
| TANAGRA | 11 | 33 | 7.0 | 53.1 | 121.6 | 73.5 |
| WEKA | 10 | 338 | 52.3 | 253.2 | 699.6 | 699.6 |

Tanagra is much faster than other programs in the decision tree induction context. This is not surprising when we know the internal data structure of Tanagra. The data are organized in columns in memory. It is very performing when we treated individually variable, when searching the discretization cut point for the split operation during the tree induction for instance.

On the other hand, in an upcoming tutorial, we see that this advantage becomes a drawback when we need to handle simultaneously all the variables during the induction process, when we perform a scalar product during the SVM (Support Vector Machine) induction for instance. In this case, Tanagra is clearly outperformed by other programs such as RAPIDMINER or WEKA.

# 4   Conclusion

This tutorial is one of the possible point views of the decision tree induction with free software in the context of large dataset. Of course, others point of views are possible: more predictors, a mixture of continuous and discrete predictors, more examples, etc. Because the programs are freely available, everyone can define and evaluate the configuration related to their preoccupation, and use their own dataset. Finally, this is what matters.

# 5   Update – October 29, 2011

One of the exciting aspects of computing is that things are changing very quickly. The machines are ever more efficient, the operating systems are improved, the software also. Since writing this document, I have a

new computer and I use a 64 bit OS (Windows 7). Some of the tools studied propose a 64 bit version (Knime, RapidMiner, R). I wonder how behave the various tools in this new context. To do that, I renewed the same experiment. We note that a more efficient computer allows to improve the computation time (about 20%). The specific gain for a 64 bit version is relatively low, but it is real (about 10%). And some tools are clearly improved their programming of the decision tree induction (Knime, RapidMiner). On the other hand, we observe that the memory occupation remains stable for the most of the tools in the new context.

So that everyone can compare the results on their own machine, here are the characteristics of my computer:
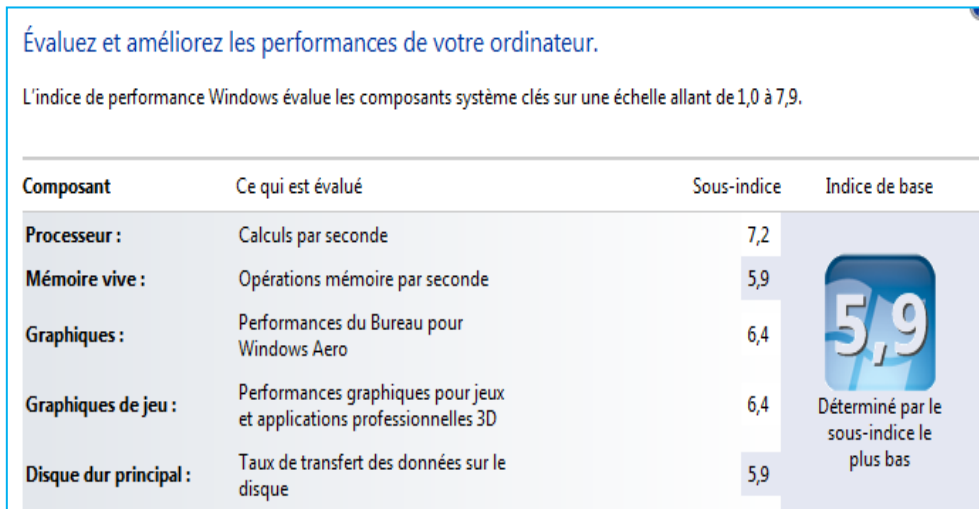
**Informations système générales**

Édition Windows

Windows 7 Édition Familiale Premium

Copyright © 2009 Microsoft Corporation. Tous droits réservés.

Service Pack 1

Obtenir plus de fonctionnalités avec une nouvelle édition de Windows 7

Système

| | |
|---|---|
| Évaluation : | **5,9** Indice de performance Windows |
| Processeur : | Intel(R) Core(TM)2 Quad CPU   Q9400  @ 2.66GHz  2.67 GHz |
| Mémoire installée (RAM) : | 4,00 Go |
| Type du système : | Système d'exploitation 64 bits |
| Stylet et fonction tactile : | La fonctionnalité de saisie tactile ou avec un stylet n'est pas disponible sur cet écran |

With the following performance index provided by Windows:

**Évaluez et améliorez les performances de votre ordinateur.**

L'indice de performance Windows évalue les composants système clés sur une échelle allant de 1,0 à 7,9.

| Composant | Ce qui est évalué | Sous-indice | Indice de base |
|---|---|---|---|
| **Processeur :** | Calculs par seconde | 7,2 | |
| **Mémoire vive :** | Opérations mémoire par seconde | 5,9 | **5,9** |
| **Graphiques :** | Performances du Bureau pour Windows Aero | 6,4 | |
| **Graphiques de jeu :** | Performances graphiques pour jeux et applications professionnelles 3D | 6,4 | Déterminé par le sous-indice le plus bas |
| **Disque dur principal :** | Taux de transfert des données sur le disque | 5,9 | |

The direct comparison is hard here because various factors influence the results. (1) We use a more powerful computer (Pentium 4 at 3.8 Ghz with 2 GB RAM vs. Core 2 Quad at 2.66 Ghz with 4 GB de RAM). (2) We have a new OS (Windows XP 32 bit vs. Windows 7 64 bit). (3) Perhaps some tools are improved their programming of the decision tree induction. So, we must remain cautious for the comments of the results.

To facilitate the reading of the results, we take the table related to our previous experiment (September 2008 – Windows XP), we put in parallel the new results (October 2011, Windows 7 - 64 bit).

| Logiciel | Tree learning processing time | | | Memory occupation (MB) | | |
|---|---|---|---|---|---|---|
| | Seconds | Rank | | Launching | Data loaded | Tree Learning |
| Knime 1.3.5 | 270 | 5 | | 93 | 160 | 246 |
| Orange 1.0b2 | 130 | 3 | | 25 | 260 | 796 |
| R 2.6.0 | 157 | 4 | | 19 | 184 | 719 |
| RapidMiner 4 | 298 | 6 | | 136 | 228 | 1274 |
| Sipina 3.0 - C4.5 | 122 | 2 | | 8 | 67 | 540 |
| Tanagra 1.4.27 | 33 | 1 | | 7 | 53 | 122 |
| Weka 3.5.6 | 338 | 7 | | 52 | 253 | 700 |

**Pentium IV 3,8 Ghz - Windows XP**

| Logiciel | Tree learning processing time | | | Memory occupation (MB) | | |
|---|---|---|---|---|---|---|
| | Seconds | Rank | Improvement (%) | Launching | Data loaded | Tree Learning |
| Knime (2.4.2 - 64 bits) - 4 threads | 39 | 2 | 86% | 131 | 328 | 393 |
| Orange (2.0b - Build 15/10/2011) | 84 | 6 | 35% | 76 | 405 | 1008 |
| R 2.13.2 (64 bits) | 79 | 5 | 50% | 28 | 192 | 719 |
| RapidMiner 5.1.011 (64 bits) | 77 | 4 | 74% | 369 | 771 | 926 |
| Sipina 3.7 - C4.5 | 52 | 3 | 57% | 3 | 63 | 76 |
| Tanagra 1.4.41 | 22 | 1 | 33% | 3 | 50 | 119 |
| Weka 3.7.4 (64 bits) | 195 | 7 | 42% | 116 | 420 | 805 |

**Q9400 2,66 Ghz - Windows 7 64 bits**

**Calculation time**. The calculation is improved whatever the tool used. We observe mainly the influence of the better efficiency of the computer. The result of Tanagra is interesting because it operates under 32 bit fashion, and I know that the source code is not improvement since the previous studied version (14.27 at September 2008, I used 1.4.41 in this new experiment). Thus, the reduction of the calculation time that we can attribute to the computer is about 30%. The result for Orange is very similar.

I think also that the source code of Weka's J48 is not really improved. Thus, we can subdivide the gain (42%) in about 30% for the computer and about 10% for the 64 bit operating system. Indeed, this new version of Weka (3.7.4) operates natively under 64 bit fashion.

To confirm this last intuition, I use the R 2.13.2 software. Two versions - 32 bit and 64 bit - are installed on the computer. From this last one, the calculation time is 79 seconds (reported into the table above). With the 32 bit version, it becomes 90 seconds. We dispose of a direct comparison here. We can conclude that the R version which operates natively in a 64 bit fashion under a 64 bit operating system improves the calculation time in 12% in relation to the 32 bit version.

**Why then such dramatic improvements for SIPINA, and especially for RAPIDMINER and KNIME?** SIPINA works under 32 bit mode. Apart the gain attributed to the computer, the improvement rests on modifications of the source code I made when I introduced a multithreaded version of the decision tree induction (see http://data-mining-tutorials.blogspot.com/2010/11/multithreading-for-decision-tree.html). I think we have the same mechanism for RapidMiner (74%) and Knime (86%). Even if they run under 64 bit mode, the improvements are really spectacular. They rest probably on a great work on the source code. We can evaluated them to about 34% [74 - (30 + 10)] for RapidMiner and 46% [86 - (30 + 10)] for Knime[1].

As regards Knime, a parameter can disturb the comparison. We use a multithreaded version in our new experiment. I set "Number threads" to 4. But, in practice, I observed that only 2 cores are really used during

---

[1] These values are approximate. They give mainly an order of magnitude on the influence of the various factors on the calculation time.

the induction process. When I set the parameter to "Number thread = 1", the calculation time becomes 48 seconds. The gain is 82%.

**About the ranking of the tools**, we observe that Tanagra remains the better. I think that this is because the internal data structure is especially well adapted to the decision tree induction. This situation is not true for another learning algorithm (e.g. SVM - http://data-mining-tutorials.blogspot.com/2009/07/implementing-svm-on-large-dataset.html). We note however that the other tools – Knime especially - are significantly improved, they are got really closer.

**Memory occupation**. We put aside SIPINA for which a great work on memory leak is made for the recent version. The fear concerns mainly the 64 bit version. The memory occupation is automatically inflated for some data structures such as pointers (see http://cran.univ-lyon1.fr/bin/windows/base/rw-FAQ.html#Should-I-run-32_002dbit-or-64_002dbit-R_003f). We observe in this new experiment that the memory occupation has not been deteriorated. On the other hand, we can theoretically handle larger datasets into main memory. For instance, for the tools which rests on the Java technology, we can set a larger value for the "-Xmx" (maximum heap size) option. The capabilities of the tools are dramatically improved.

**Conclusion**. The main conclusion of this update is that, on a 64 bit system, for the same computer configuration, a native 64 bit application is about 10% faster in a decision tree induction process compared to the same program compiled in a 32 bit mode[2]. The advantage of the 64 bit application is all the more interesting that it can theoretically handle large datasets into main memory. And for the same dataset, the memory occupation is marginally increased. The second conclusion is that some tools, mainly Knime and RapidMiner, have considerably improved their program for the decision tree induction over the 3 last years. But, to date, and this is the last conclusion, Tanagra nevertheless remains very fast for the induction of decision tree.

---

[2] We note however that this conclusion is not definitive. Our dataset has some particular characteristics: all the descriptors are continuous, there are a few number of descriptors and a large number of instances. We should extend the experiment to other kind of dataset to obtain a larger scope results.