# 1. Introduction

**Working with Sipina under Linux using Wine**.

In a recent tutorial ([http://data-mining-tutorials.blogspot.com/2009/01/tanagra-under-linux.html](http://data-mining-tutorials.blogspot.com/2009/01/tanagra-under-linux.html)), we show that it is possible to work with Tanagra under Linux using Wine. In this document, we implement Sipina with the same framework i.e. we install and use Sipina in a Linux environment. We use the Ubuntu distribution (French version 8.10). All the functionalities of Sipina are available, especially the interactive tools which allows us to explore deeply the subpopulation into a node of the tree.

In this tutorial, we implement the following steps:

1. Installing Sipina under Linux;
2. Launching the software;
3. Loading a dataset (text file with tab separator);
4. Choosing the class attribute and the predictive variables;
5. Partitioning the dataset in a train set and test set;
6. Computing the tree on the train set;
7. Evaluation the tree on the test set e.g. computing the confusion matrix, the error rate, etc.;
8. Exploring a subpopulation related to a node of the tree;
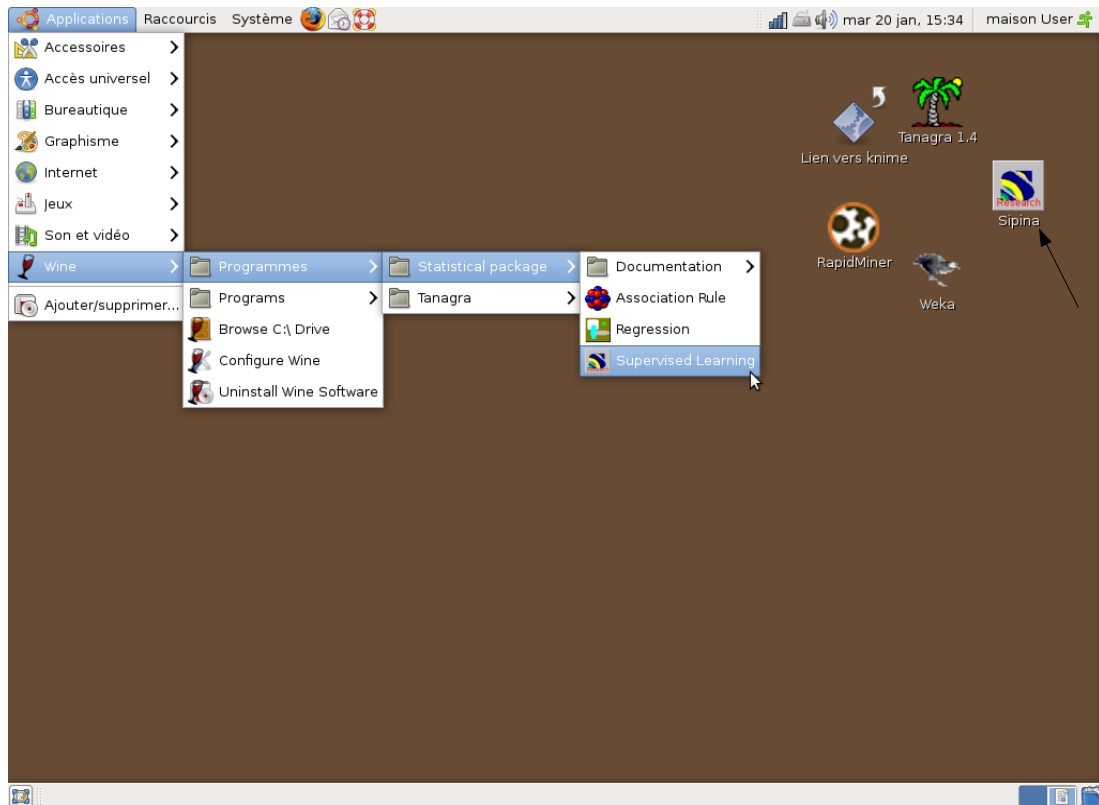9. Launching a new analysis on a subpopulation related to a node of the tree.

We will describe quickly the various features of the software in this tutorial. They are already presented in several documents available online ([http://eric.univ-lyon2.fr/~ricco/sipina.html](http://eric.univ-lyon2.fr/~ricco/sipina.html), see the DOWNLOAD section). Our main goal here is to show the capabilities of Sipina under Linux.

We use the French Ubuntu 8.10 distribution ([http://www.ubuntu.com/](http://www.ubuntu.com/)); we have installed also Wine [[http://en.wikipedia.org/wiki/Wine_(software)](http://en.wikipedia.org/wiki/Wine_(software))]], a program which allows to Windows programs to run under Linux. We show in the tutorial that Sipina is fully functional in this environment.
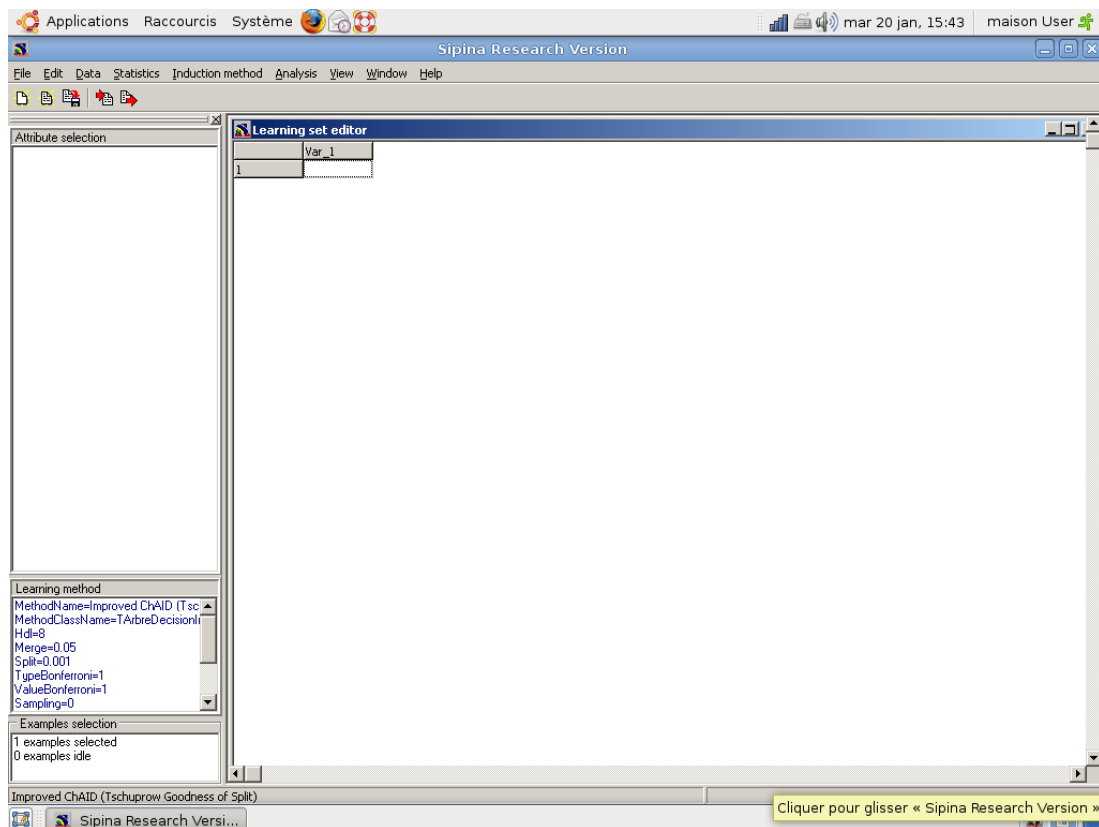
# 2. Installation of Sipina

After we download the setup file « setup_statp_ackage.exe » on our server ([http://eric.univ-lyon2.fr/~ricco/sipina.html](http://eric.univ-lyon2.fr/~ricco/sipina.html)), we start the installation using wine. The simplest way is to enter the following command line on a terminal window: *wine /your path/setup_stat_package.exe*

The installation process is started. We can validate each step and use the default installation path (c:\program files\...). At the end of the installation, the "Statistical Package" is now available in the "Applications" menu as we can see in the screenshot below. Sipina corresponds to the "Supervised Learning" item. We can also create a shortcut on the desktop.
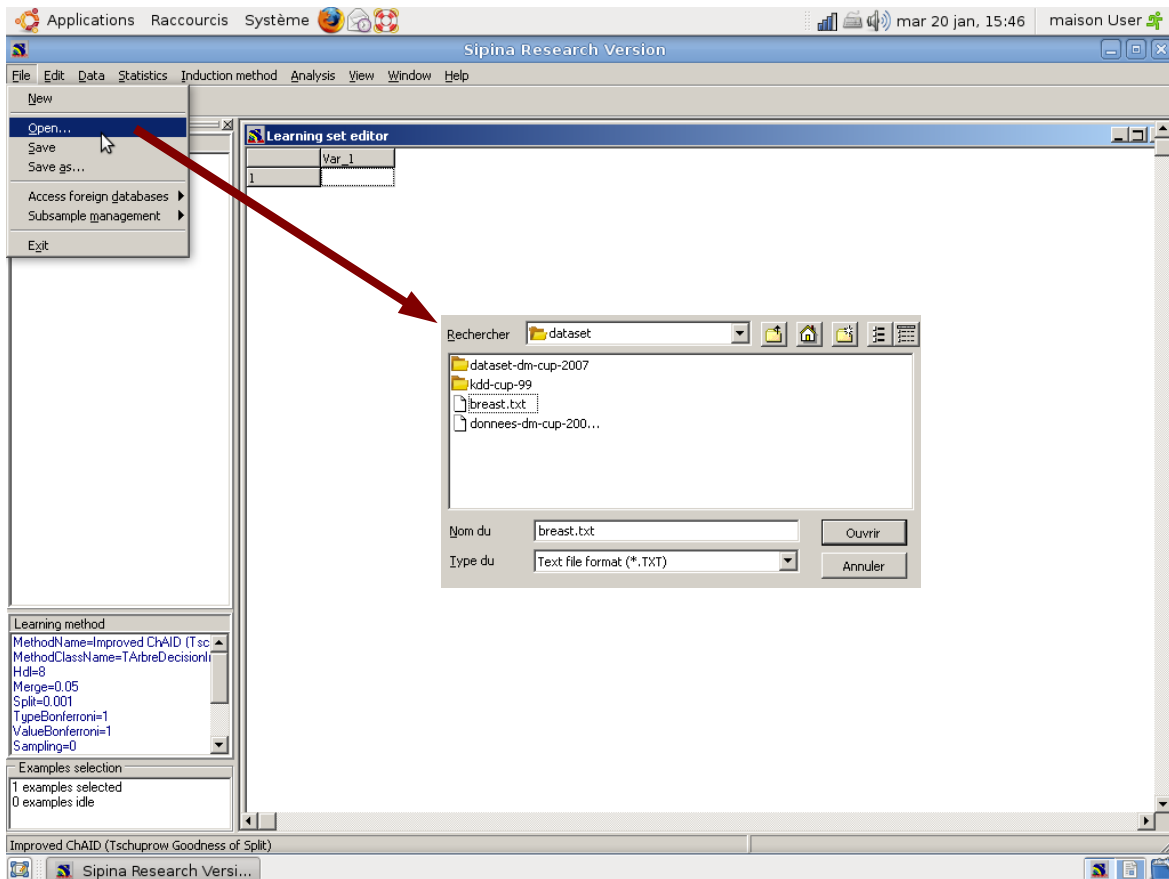
# 3. Using Sipina under Linux



When we start Sipina, we find the usual interface. Nothing distinguishes it from its running under Windows.

**Data importation**. We use the BREAST.TXT in this tutorial (http://eric.univ-lyon2.fr/~ricco/tanagra/fichiers/breast.txt ; *Note: the decimal separator is always ". " for Sipina, whatever your system configuration*). We want to predict the characteristics of cells extracted from a tumor ("malignant" or "benign") from their description (shape, size, etc.). To import the dataset, we click on the FILE/OPEN menu. We select the text file format.



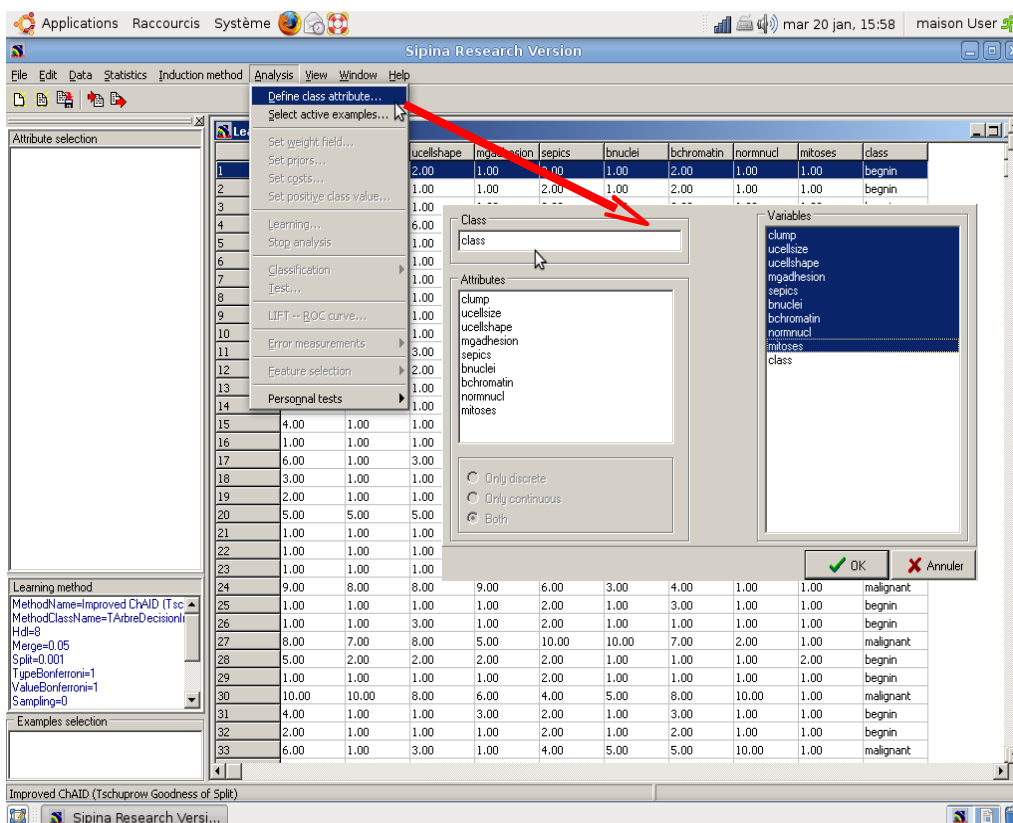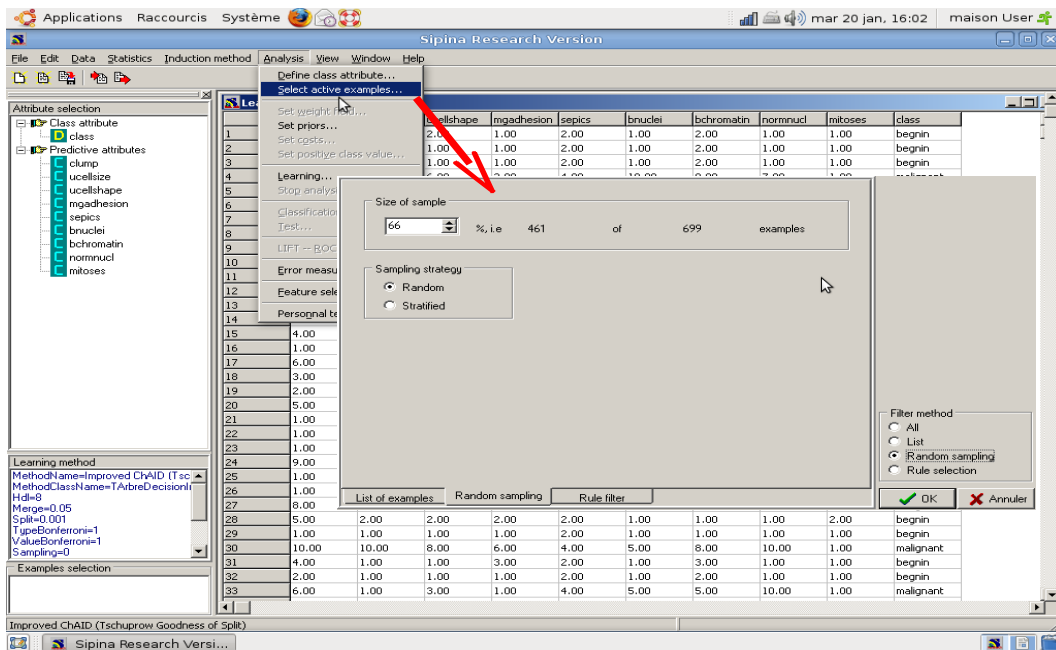In the dialog box which appears, we set: the column separator and the variables name (in the first row).

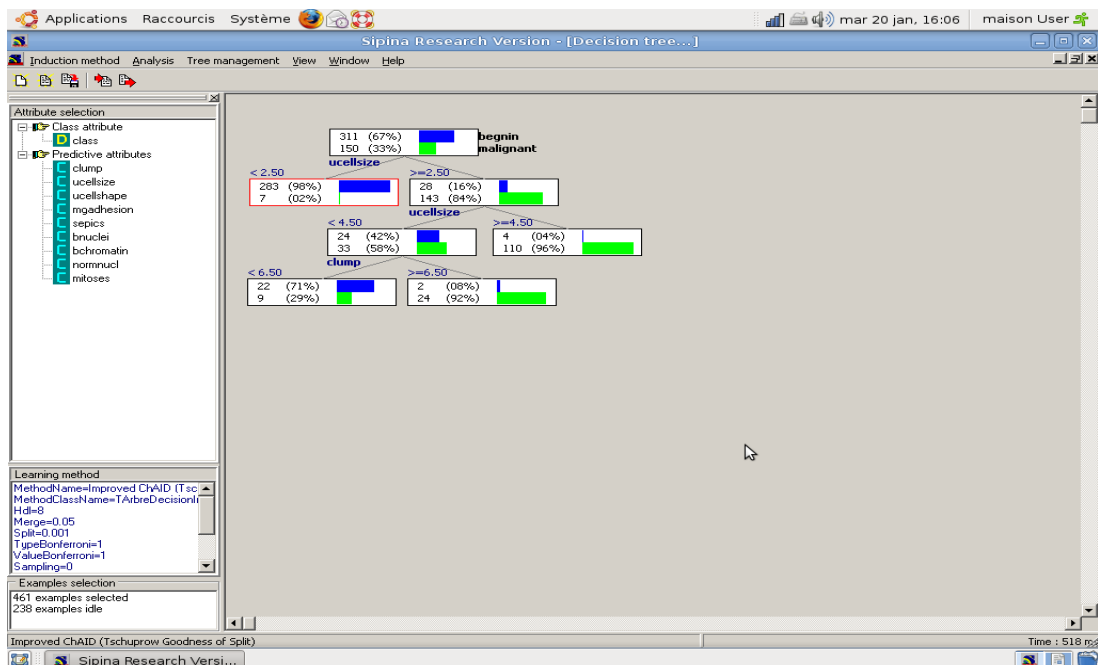The dataset is displayed in a grid.



**Defining the status of the variables.** In order to define the class attribute, we click on the ANALYSIS / DEFINE CLASS ATTRIBUTE menu. We use drag-and-drop to select the variables.

**Partitioning the dataset**. We want to use 66% of the dataset for the learning phase, the remainder (34%) for the test phase. We click on the ANALYSIS / SELECT ACTIVE EXAMPLES menu. Into the dialog box, we choose the RANDOM SAMPLING item: 66% of the dataset is selected i.e. 461 examples.
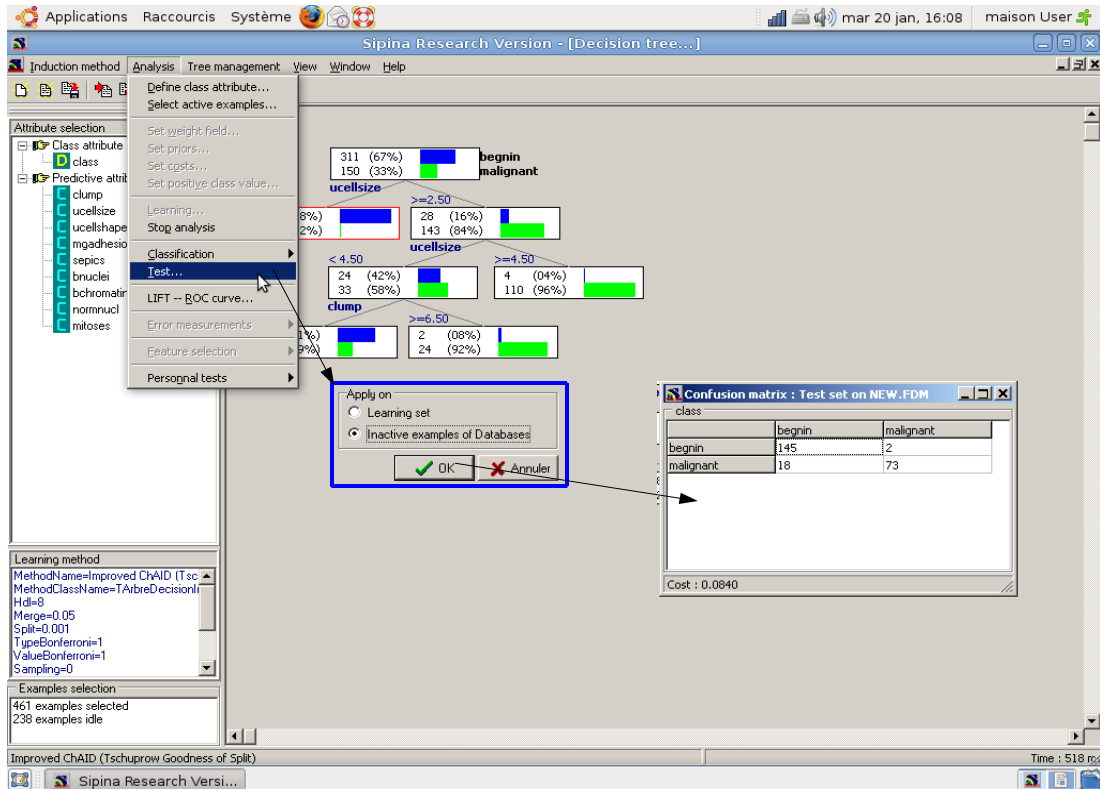


**Creating the decision tree.** We click on the ANALYSIS / LEARNING menu to start the analysis. The tree is displayed in a new window.
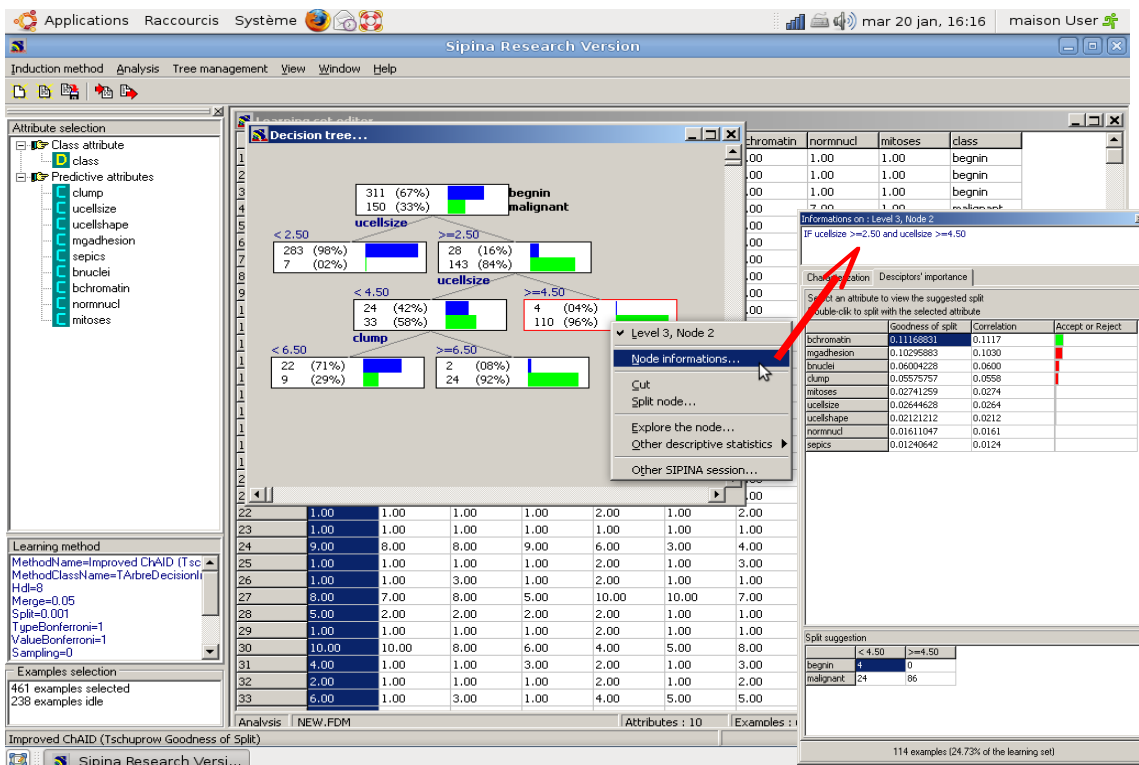


*Note: Because the data is randomly partitioned, perhaps we do not obtain the same tree.*

**Evaluation on the test set**. We want to evaluate the classifier performance on the remainder dataset i.e. 238 examples. We click on the ANALYSIS / TEST menu. In the dialog box, we select the "Inactive

Examples". The confusion matrix is displayed in a new window. The test error rate is 8.40%.



**Exploration of a subpopulation**. We want to highlight the characteristics of a subpopulation related to a node. We select the node; we click on the NODE INFORMATIONS item into the contextual menu. In the DESCRIPTORS IMPORTANCE tab, we have the goodness of split of each predictive variable.

Into the CHARACTERIZATION tab, we can compare the descriptive statistics computed on the whole population (the root of the tree) and the studied subpopulation (the current node of the tree).

For instance, we observe that UCELLSHAPE, which is not used in the tree, allows characterizing differently the current node: the local mean (7.57) of this variable is higher here than the global mean (3.16) computed on the whole dataset.

We observe that all the variable have significantly higher values into this subpopulation.



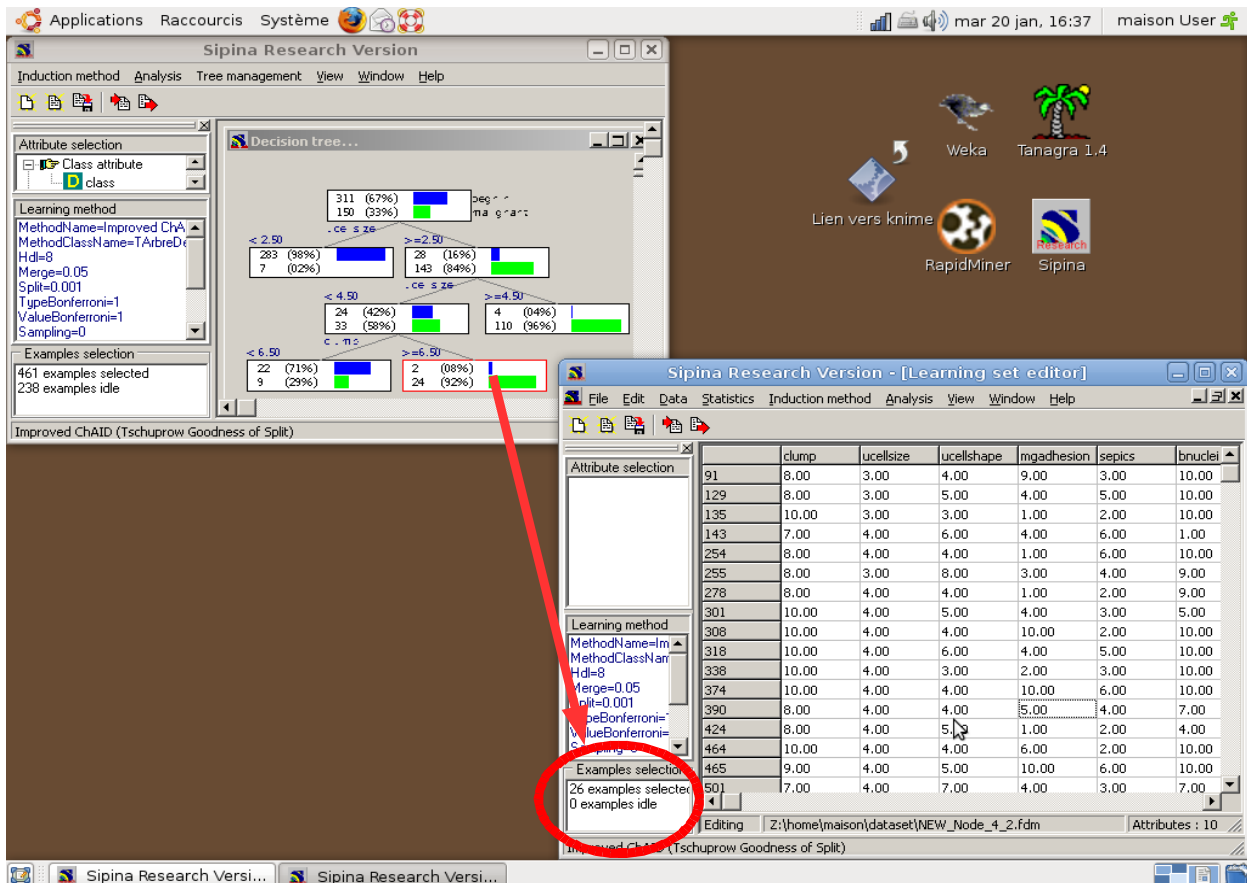## 4. Exploring more deeply a subpopulation

Sometimes, we want to explore mode deeply the node by launching a new analysis on the related subset of examples. For instance, we want to study the subset of examples (26 cases) described as follows:

« UCELLSIZE >= 2.5 **and** UCELLSIZE < 4.5 **and** CLUMP >= 6.5 »

With the contextual menu, we click on the OTHER SIPINA SESSION item. A new SIPINA session is started and the corresponding subset of examples is automatically loaded.

All the variables, but only the selected examples, are sent. We can start now a new analysis.

# 5. Conclusion

In this tutorial, we showed that it is possible to install and use Sipina in the Linux environment without the need for specific technical skill. We then have the whole environment of the software which is fully functional.