# 1   Theme

**Regression model deployment using Tanagra.**

Model deployment is one of the main objectives of the data mining process. We want to apply a model learned on a training set on unseen cases i.e. any people coming from the population. In the classification framework, the aim is to assign to the instance its class value from their description (e.g.   http://data-mining-tutorials.blogspot.com/2008/11/apply-classifier-on-new-dataset.html).   In the clustering framework, we try to detect the group which is as similar as possible to the instance according their characteristics (e.g. http://data-mining-tutorials.blogspot.com/2008/12/k-means-classification-of-new-instance.html).

We are concerned about the regression framework here[1]. The aim is to predict the values of the dependent variable for unseen instances (or unlabeled instances) from the observed values on the independent variables. The process is rather basic if we handle a linear regression model. We apply the computed parameters on the unseen instances. But, it becomes difficult when we want to treat more complex models such as support vector regression with nonlinear kernels, or the models elaborated from a combination of techniques (e.g. regression from the factors of a principal component analysis). In this context, it is essential that the deployment process is directly ensured by the data mining tool.

With Tanagra, it is possible to easily deploy the regression models, even when they are the result of a combination of technique. Simply, we must prepare the data file in a particular way. In this tutorial, we describe below how to organize the data file in order to deploy various models in an unified framework: a linear regression model, a PLS regression model, a support vector regression model with a RBF (radial basis function) kernel, a regression tree model , a regression model from the factors of a principal component analysis. Then, we export the results (the predicted values for the dependent variable) in a new data file. Last, we check if the predicted values are similar according to the various models.
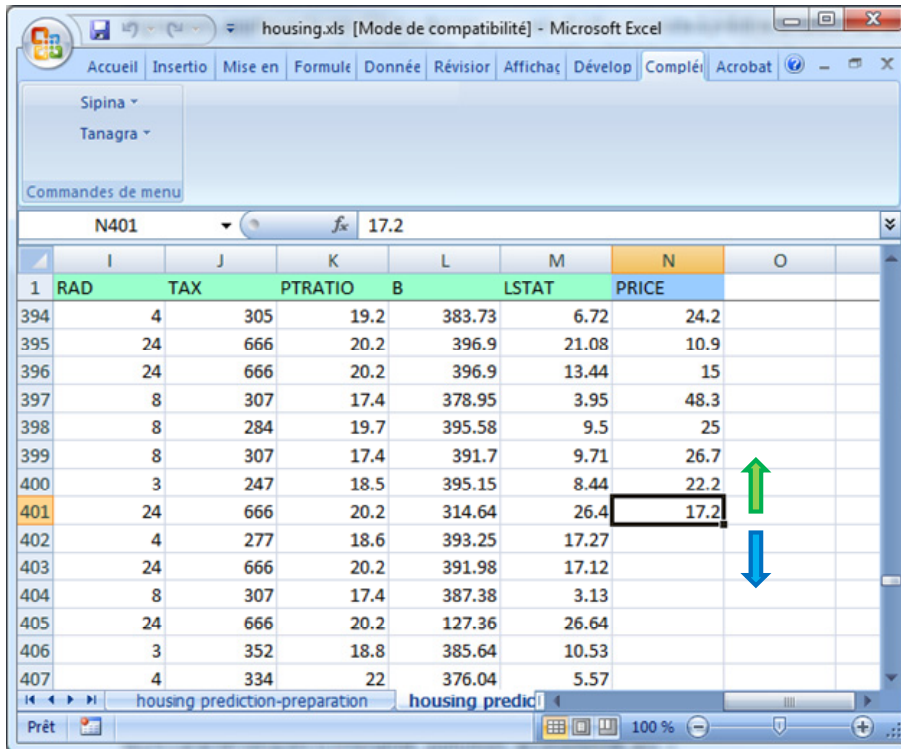
# 2   Organizing the data file

We use the HOUSING dataset (http://archive.ics.uci.edu/ml/datasets/Housing). We want to predict the values of PRICE from the house characteristics (criminality, pollution, etc.). We have 400 instances (the training sample) for the construction of the models. Then, we want to apply them to 160 unseen instances i.e. the instances for which we have only the values of the independent variables.
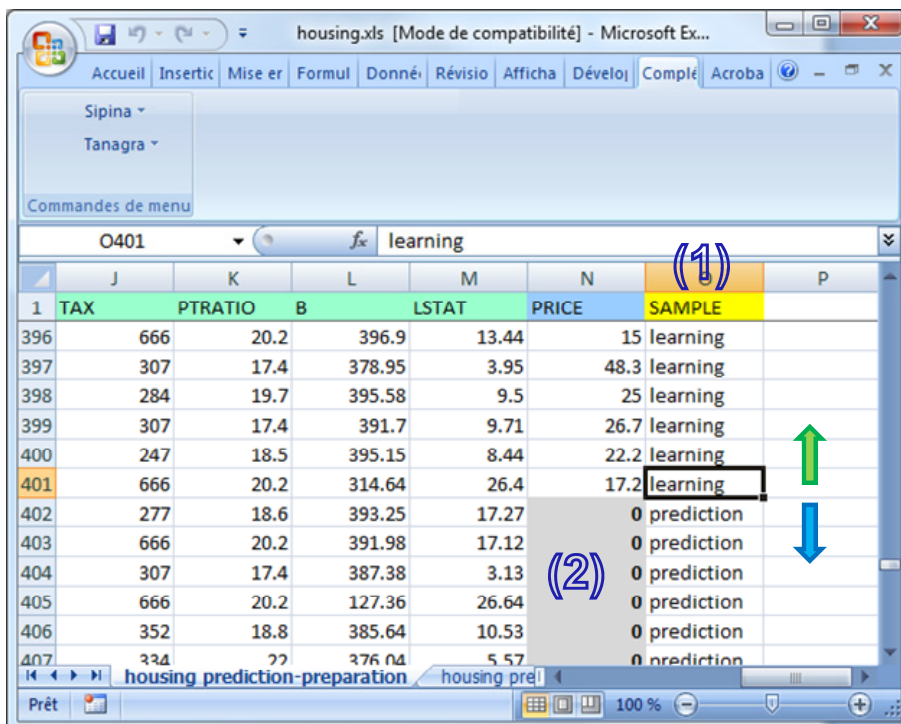
A natural way to organize the dataset is to put the training sample in the first rows of the data file. Then, we set the 106 unlabeled cases in the following rows. But, because Tanagra does not handle missing values, this data file will be truncated during the importation process. These unseen instances are not imported. Therefore, we cannot apply the models to these instances. And yet, this is our main goal.

---

[1] http://data-mining-tutorials.blogspot.com/search/label/Regression%20analysis

So that Tanagra can handle the data correctly, we need to change the organization of the data file. First, we add the SAMPLE column. It specifies the status of observations: "learning" corresponds to the learning set; "prediction" corresponds to the unseen instants. Second, we set a default value for the dependent variable. The value o (zero) is entirely appropriate for that. The aim is to circumvent the missing value problem.
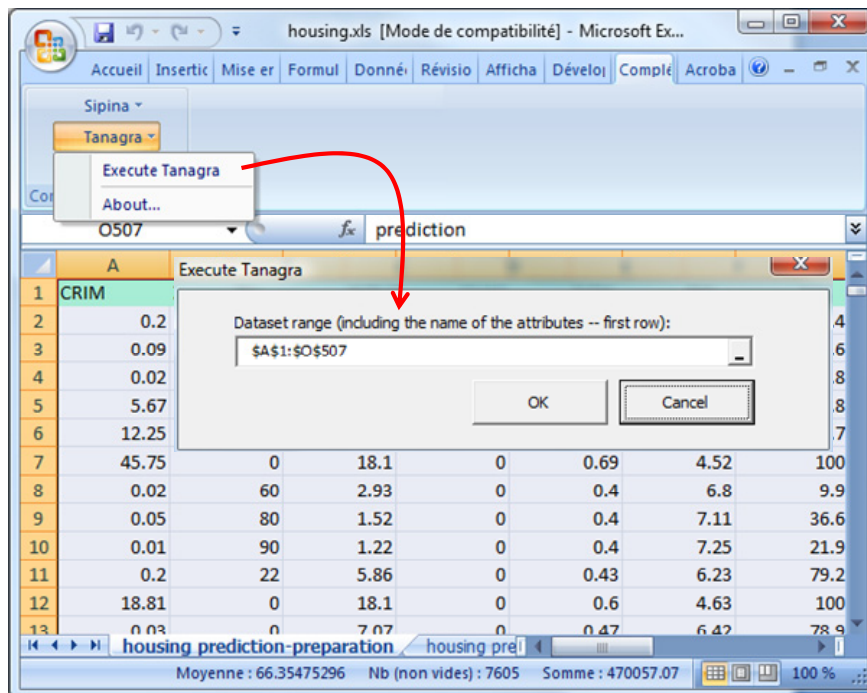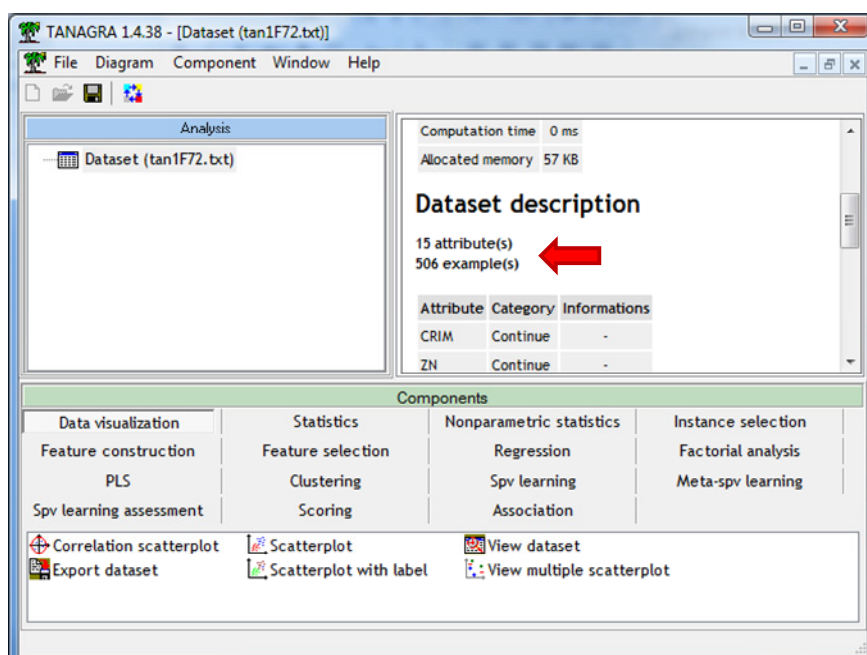


Thus, we can import the dataset into Tanagra.

# 3 Regression model deployment

## 3.1 Importing the data file

We load it into the Excel spreadsheet (we can also use OpenOffice or LibreOffice[2]). Then, using the Tanagra.xla add-in, we send the dataset to Tanagra (*see* *http://data-mining-tutorials.blogspot.com/2010/08/tanagra-add-in-for-office-2007-and.html* *for Excel 2007 and 2010;* *http://data-mining-tutorials.blogspot.com/2008/10/excel-file-handling-using-add-in.html* *for up to Excel 2003 version*).
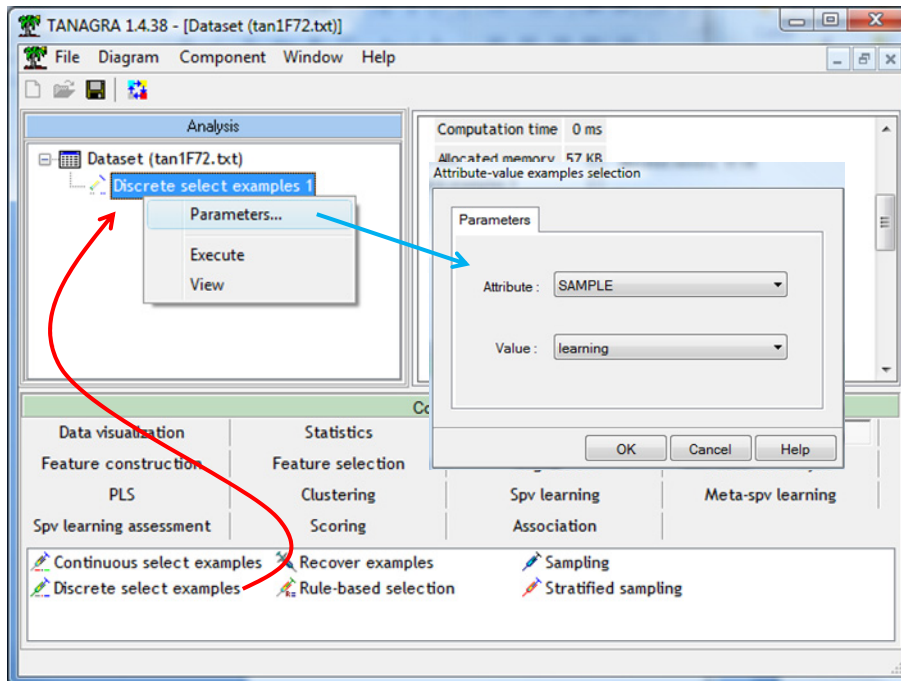
All the instances (400 + 106 = 506) are now available into Tanagra.

## 3.2    Partitioning the data file

We insert the DISCRETE SELECT EXAMPLES component (INSTANCE SELECTION tab) component into the diagram. It allows to specify the training sample i.e. the instances used for the construction of the regression model. The selection attribute is the SAMPLE column. LEARNING allows to specify the training sample.



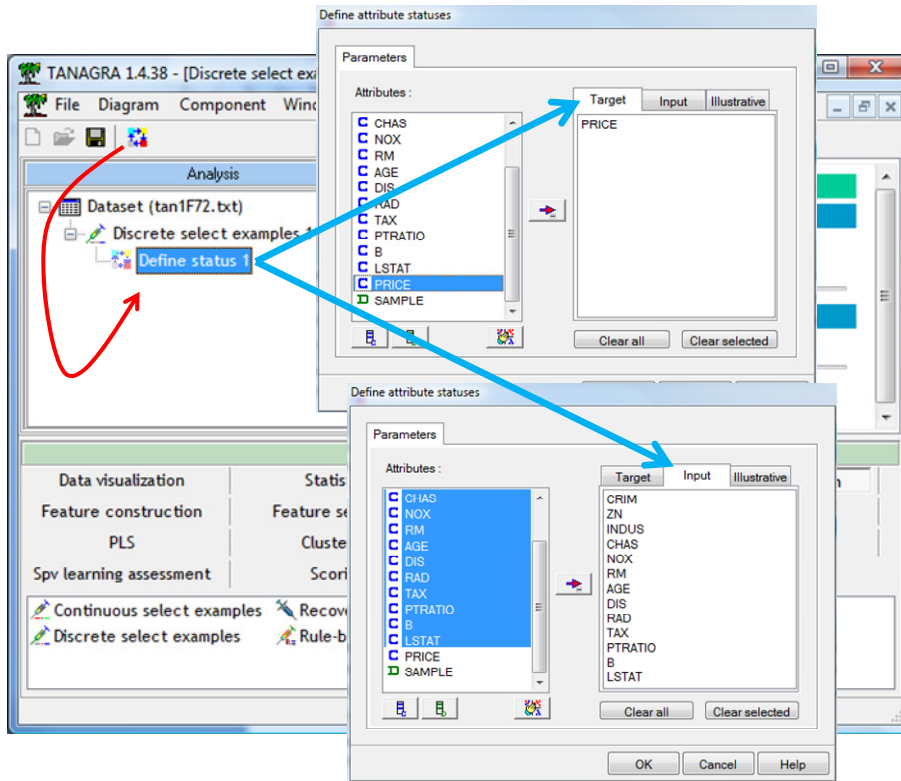Now, 400 instances are used for the learning phase.



The remaining rows are not used during learning. However, and this is the secret of the process, when the component applies the model on the data, it does on all the rows, including the unselected observations. **We take advantage to this property to get predictions on the unlabeled instances**.
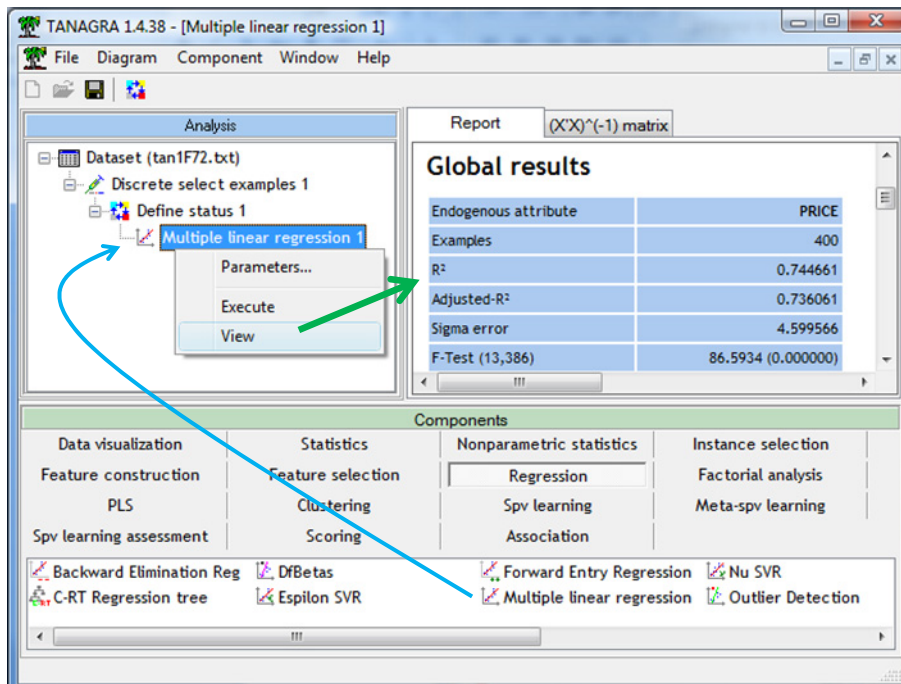
## 3.3    Construction of the models

### 3.3.1    Linear regression

We insert the DEFINE STATUS into the diagram in order to define the role of the variables into the learning process: PRICE is the TARGET attribute; the others are the INPUT ones. The SAMPLE column is no longer used at this step.
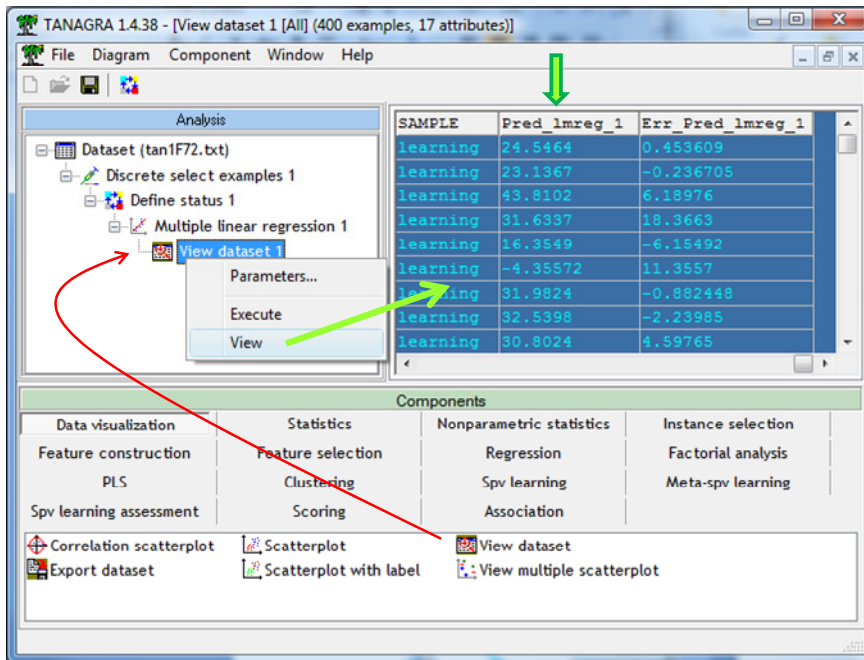


Then, we add the MULTIPLE LINEAR REGRESSION component (REGRESSION tab). We click on the VIEW menu to obtain the regression results.
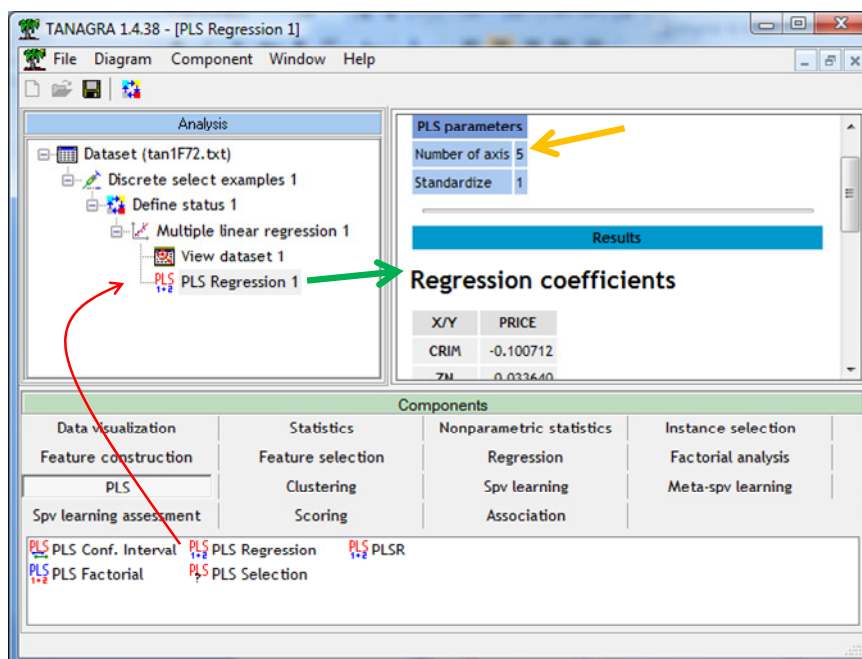
The coefficient of determination is R2 = 0.744. Some independent variables are not significant at 5% significance level. We will not detail the analysis of the results in this tutorial. Our goal is to describe the deployment of learned models on unlabeled instances.

We can visualize the dataset using the VIEW DATASET (DATA VISUALIZATION tab) component. Two columns are added. PRED_LMREG_1 is the prediction of the model. Only the prediction for the learning set is displayed for the moment. But the same calculations are performed on the unselected instances. ERR_PRED_LMREG_1 corresponds to the residuals.
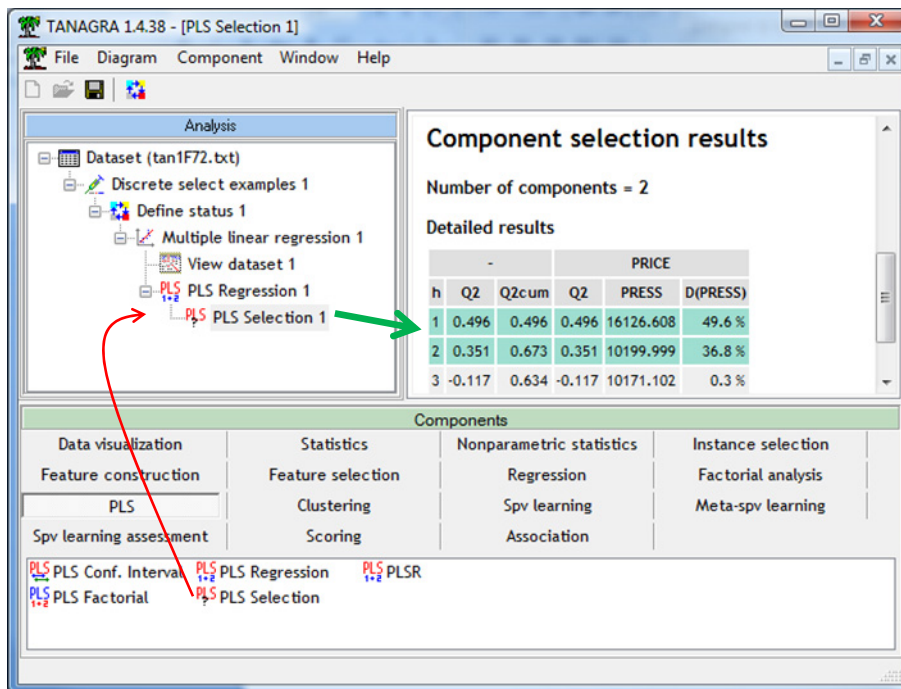


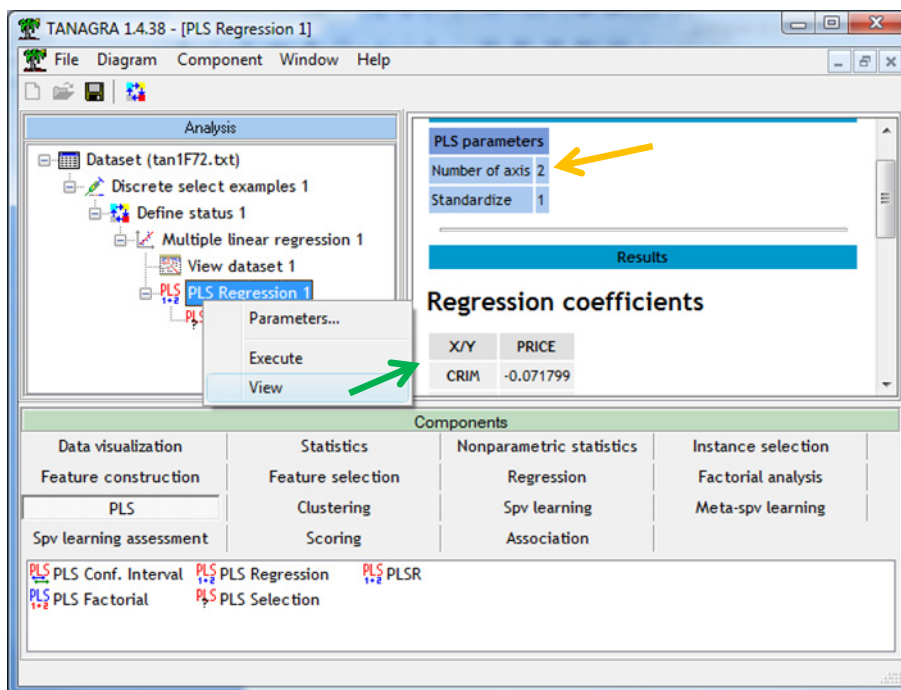### 3.3.2    PLS (Partial Least Squares) Regression

We add the PLS REGRESSION component (PLS tab) behind the linear regression. We click on the VIEW menu to obtain the results according to the default settings.
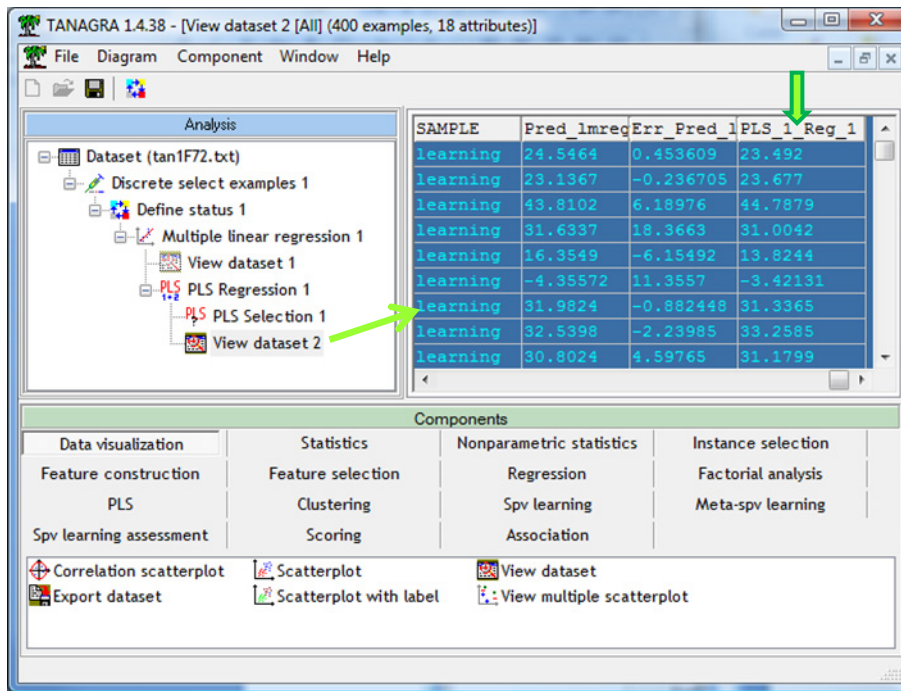
By default, 5 factors are computed. This choice is not necessarily appropriated. We can ask to Tanagra to detect the right number of factors using a cross-validation resampling approach. We use the PLS SELECTION component for that.
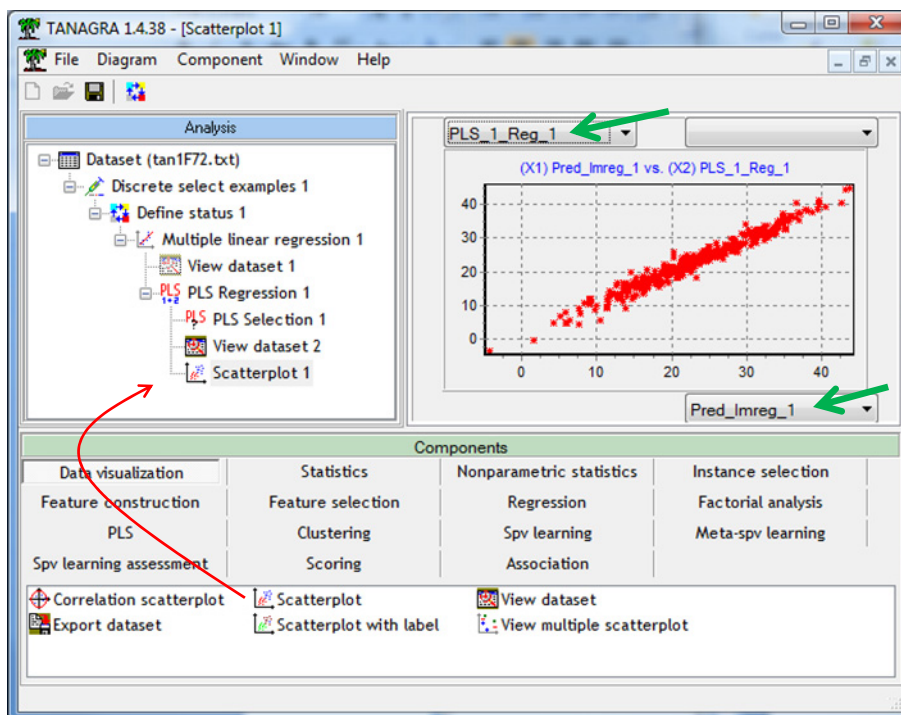


2 axes seem to be a good compromise. When we click on the VIEW menu of the PLS REGRESSION component, we obtain the definitive model coefficients.



Here also, when we visualize the dataset, we observe that the prediction column **PLS_1_REG_1** is added.

We can compare the predictions of the two models (Linear Regression and PLS Regression) on the learning set using the SCATTERPLOT (DATA VISUALIZATION tab) component.
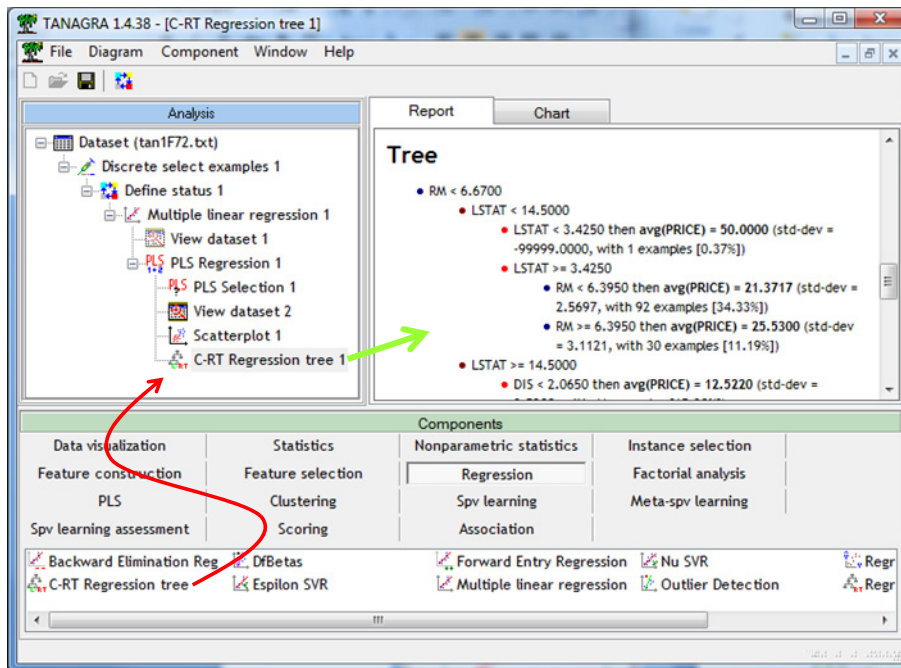


We observe that they are consistent. This is not really surprising. These are both linear models. Their behavior is quite similar on the majority of problems. The PLS1 regression seems especially valuable when we have correlated independent variables (e.g. http://data-mining-tutorials.blogspot.com/2010/05/solutions-for-multicollinearity-in.html).

### 3.3.3    Regression tree

Now, we add the regression tree component (CART approach) into the diagram. We use the C-RT Regression Tree (REGRESSION tab).
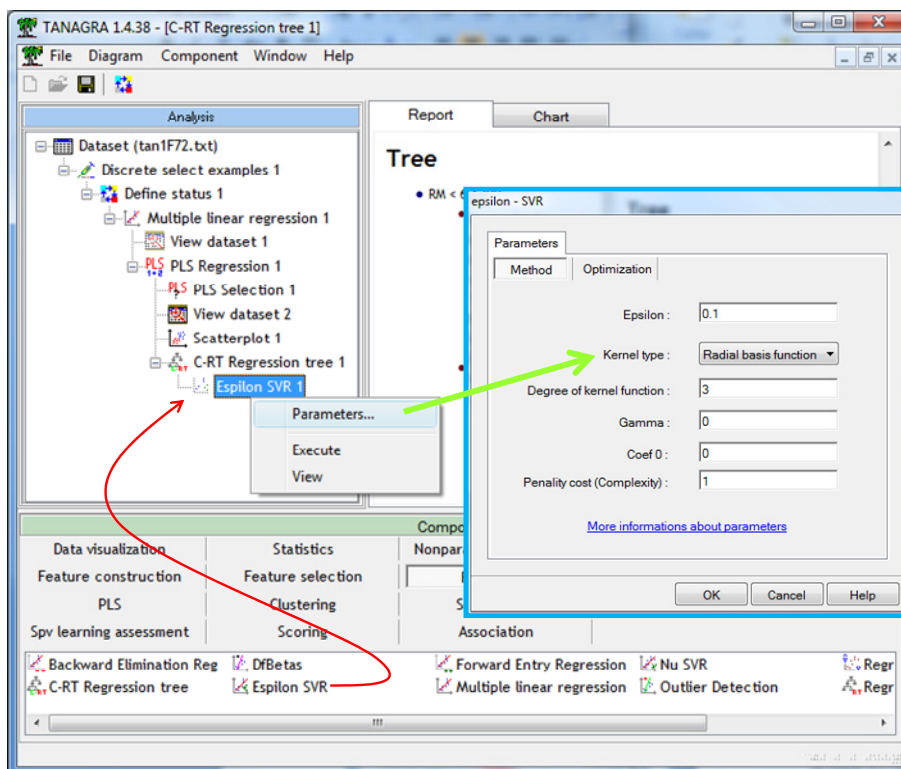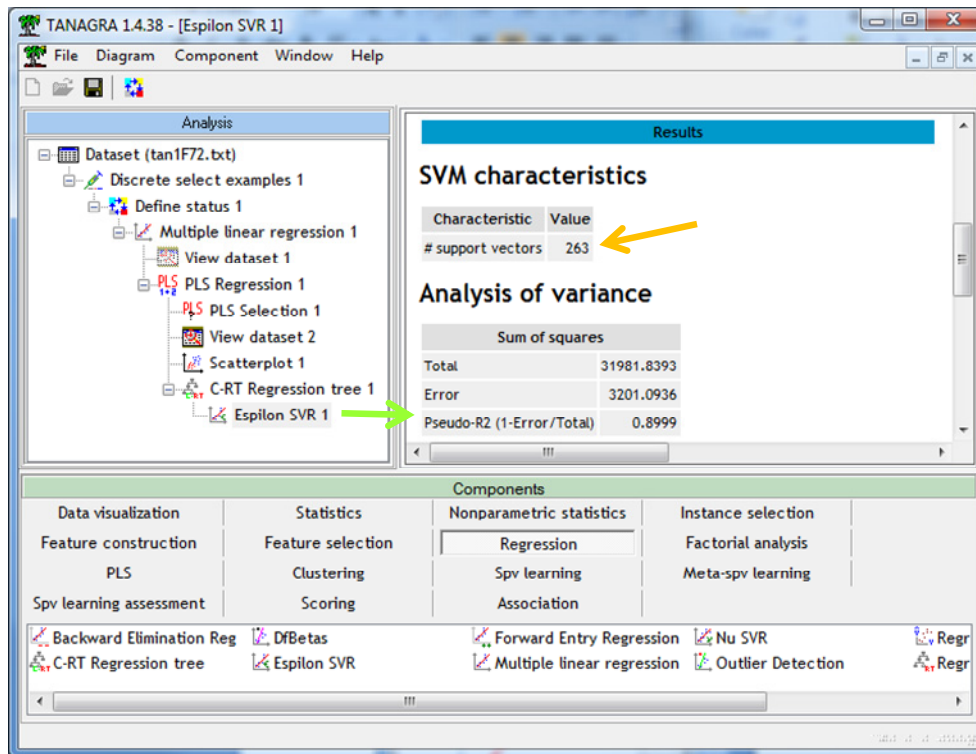
The tree consists of 12 leaves. Like the linear models, the deployment of a regression tree is very easy. We can transform the tree in a set of rules without loss of information.

### 3.3.4    Support vector regression

We use the EPSILON SVR component (REGRESSION tab). We have described the behavior of this component elsewhere (this component is based on the LIBSVM library - http://data-mining-tutorials.blogspot.com/2009/04/support-vector-regression-svr.html). We set a RBF (Radial Basis Function) kernel.
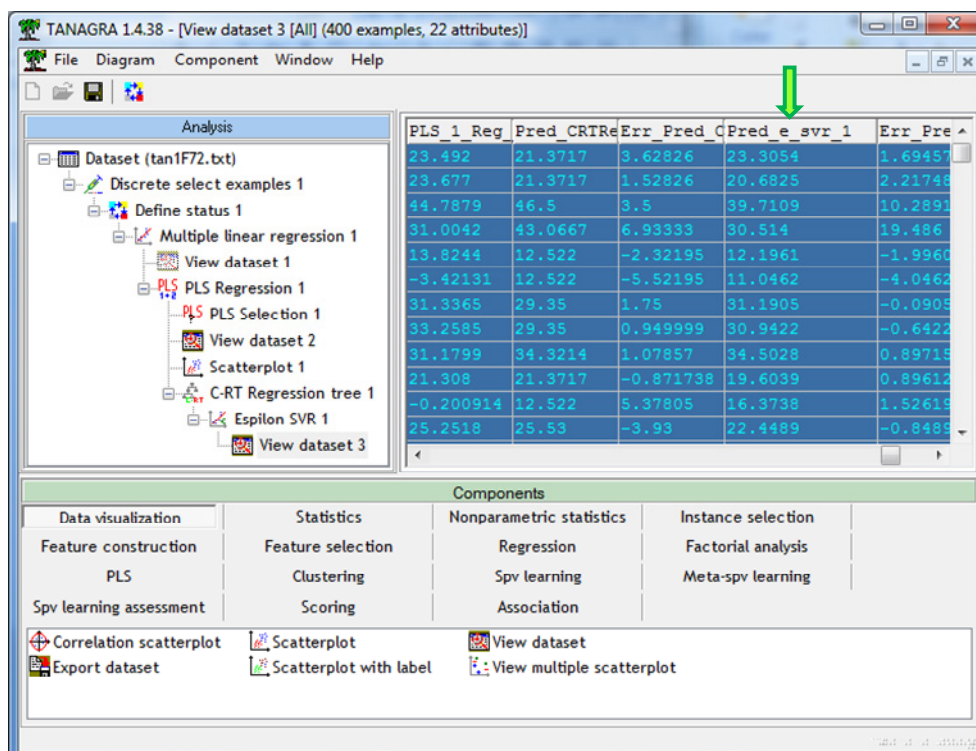


We click on the VIEW menu to obtain the results.

Here, the difficulties really begin. Indeed, we do not have an explicit model. We need the support vectors (263 instances for our problem according the output of the tool) to make predictions on unlabeled instances. It is not really easy to make this outside the data mining tool used for the learning of the model. This is the reason for which Tanagra provides functionality for the prediction on unlabeled instances. Almost all the well known data mining tools perform this kind of operation.
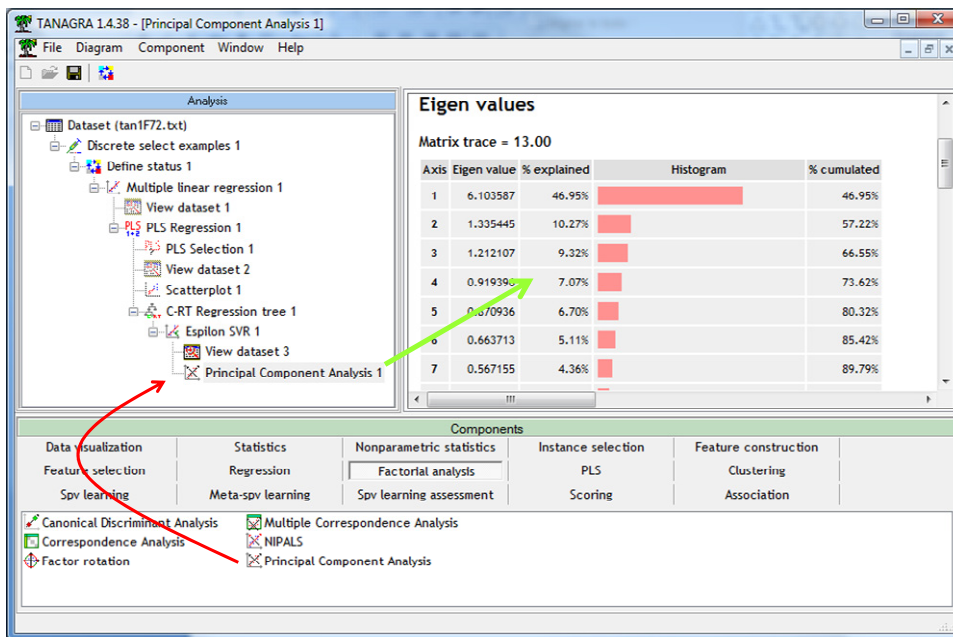
With the VIEW DATASET component, we observe that EPSILON SVR adds also new columns for the predictions and the residuals.
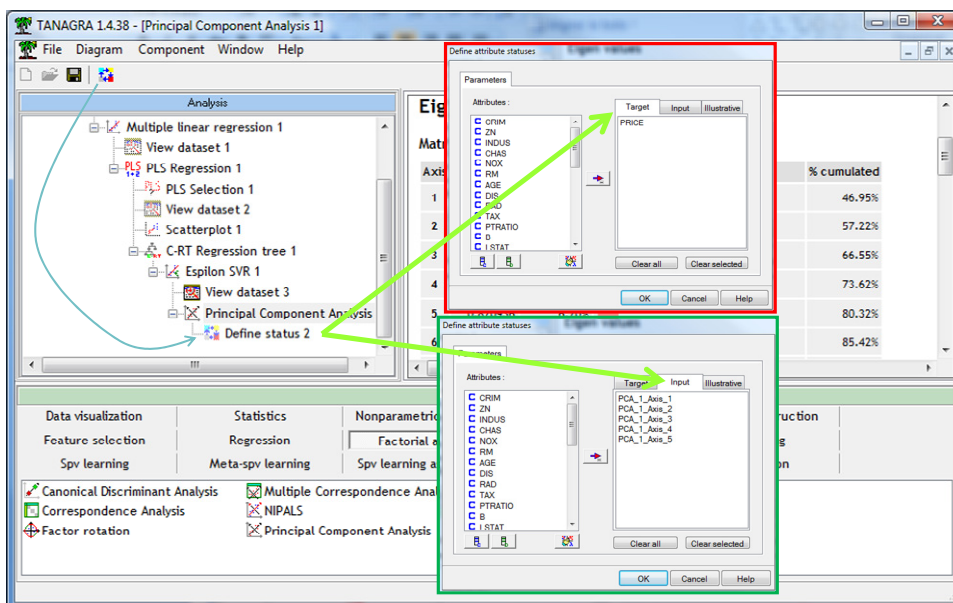
### 3.3.5    Regression on the factors of PCA

Until now, we use the learning algorithms individually. The deployment involves a single model. In this section, we want to implement a combination of methods. In a first step, we perform a principal component analysis from the independent variables. In a second step, we use the 5 first factors scores provided by the PCA as input variable in a linear regression approach. The dependent variable is still PRICE of course. For the deployment phase, we must sequentially use the two models. First, we compute the factor scores of unlabeled instances. Then, we obtain the prediction by applying the regression parameters of these factor scores.
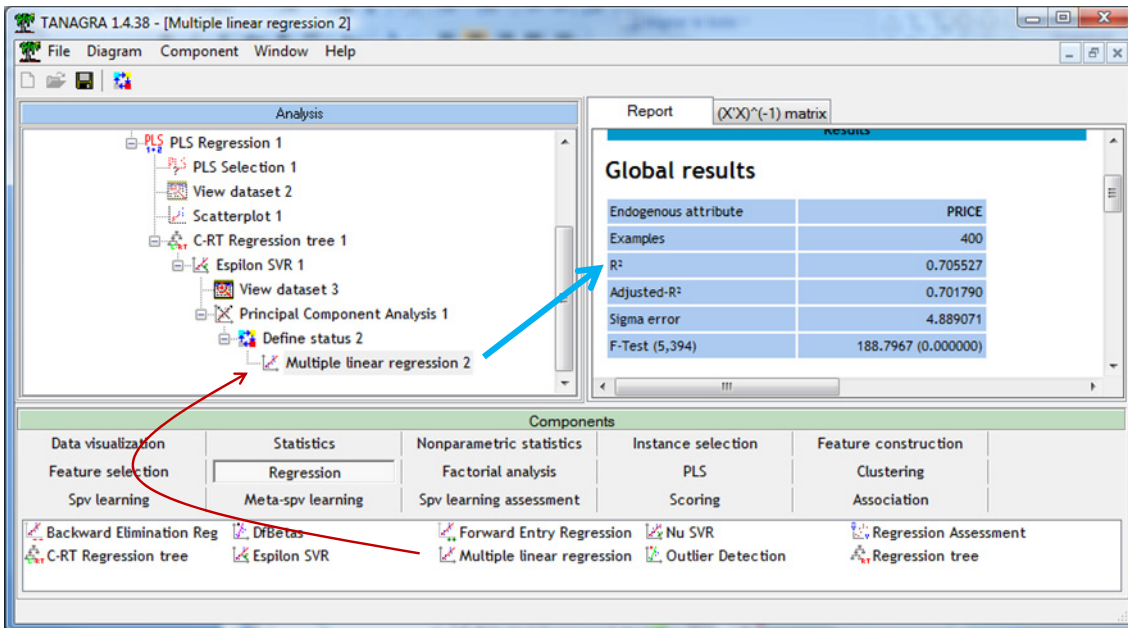
We add the PRINCIPAL COMPONENT ANALYSIS component (FACTORIAL ANALYSIS tab) into the diagram. We obtain the following results.



For the regression step, with the DEFINE STATUS component, we set the 5 first factors as INPUT variables, PRICE is the TARGET one.
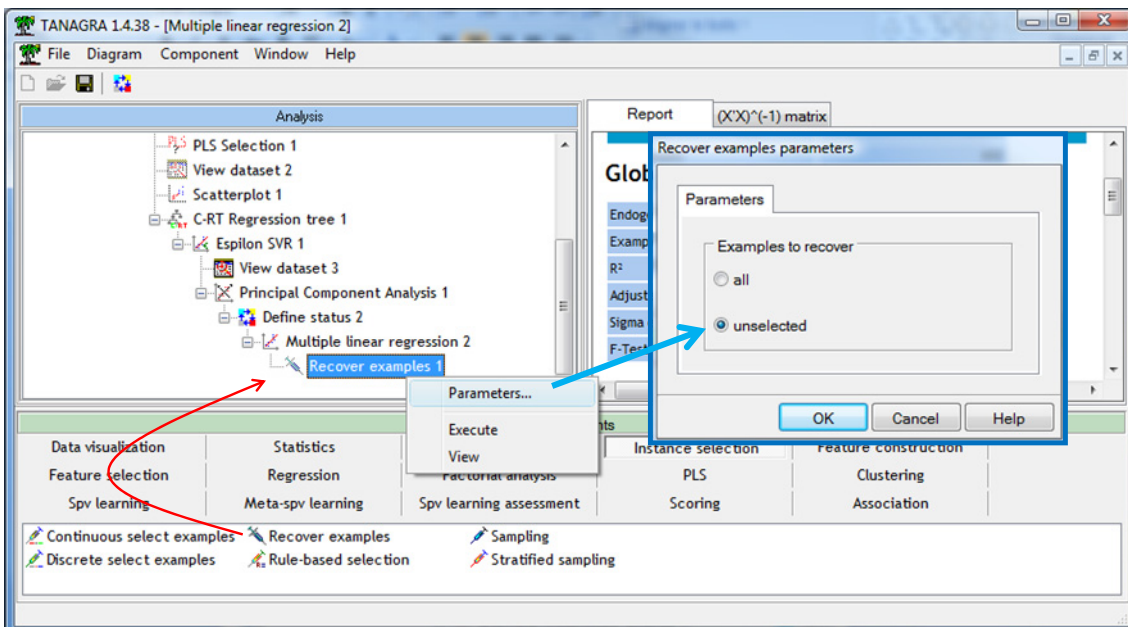
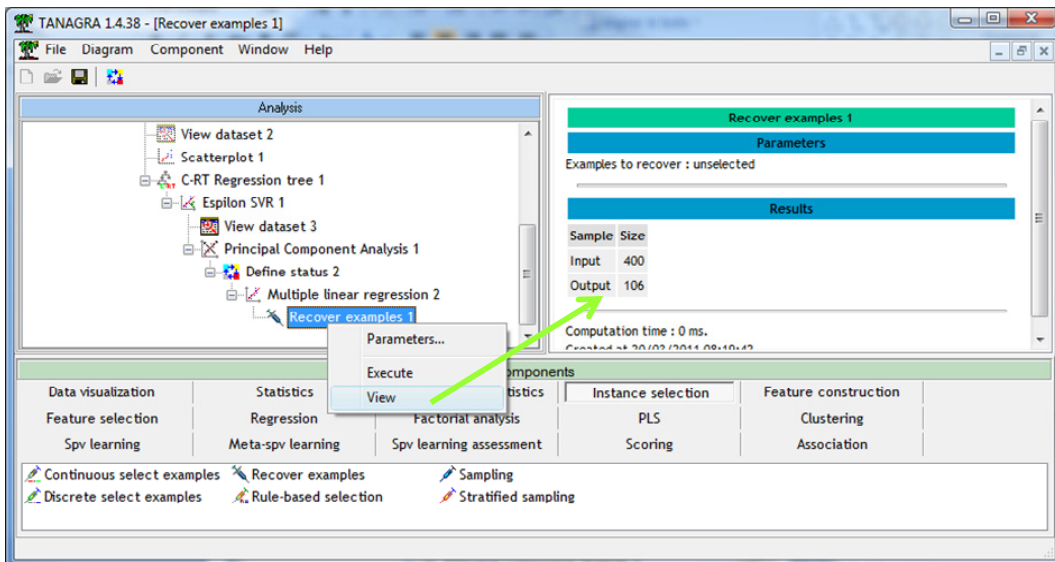Then, we insert the linear regression component. We click on the VIEW menu.



### 3.4 Retrieving the predictions for the unlabeled instances

For each regression component, TANAGRA performs the predictions on the unselected instances i.e. on the unlabeled instances according our data file organization. To visualize them, we must reverse the instances selection.
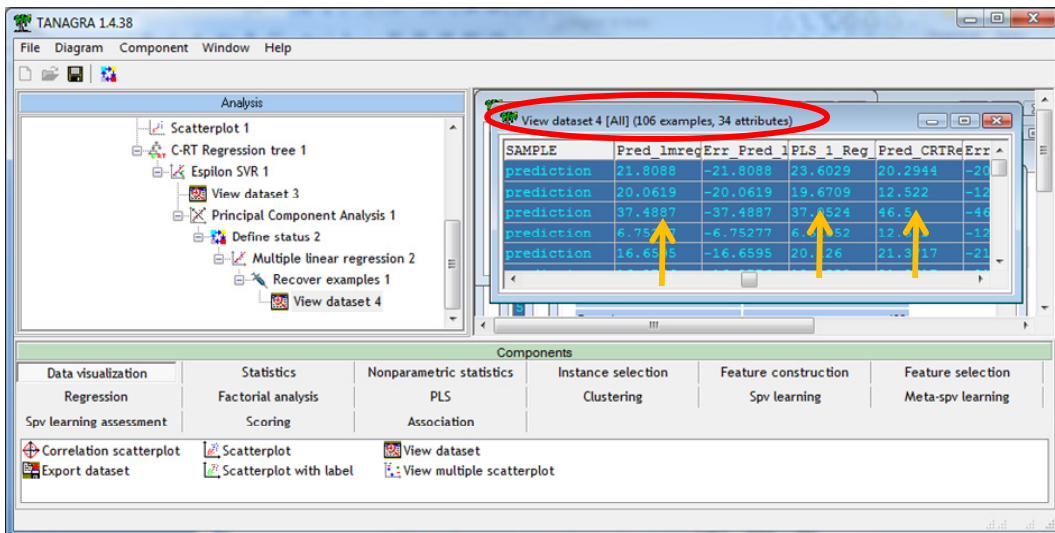
We insert the RECOVER EXAMPLES (INSTANCES SELECTION tab) into the diagram. We set the following settings.
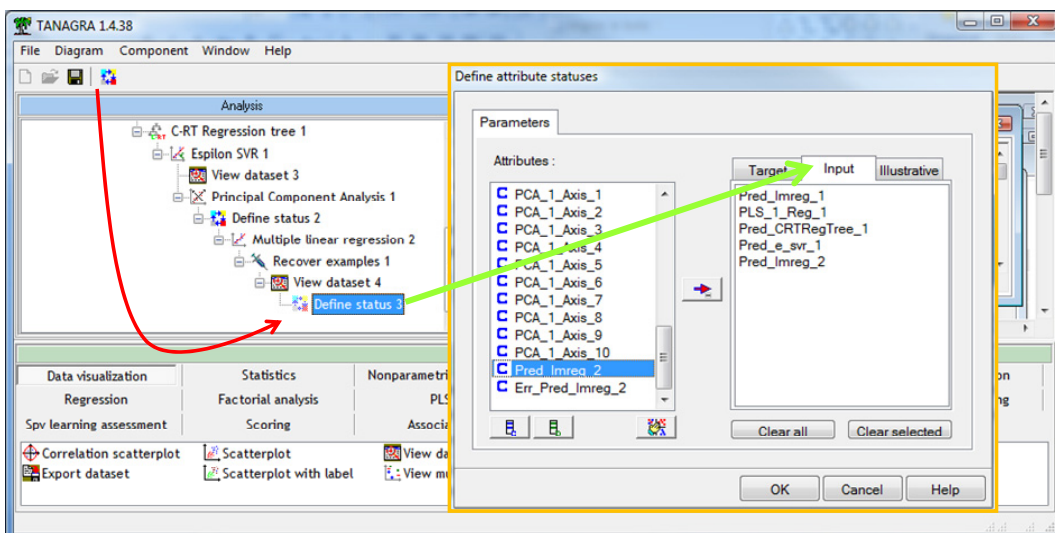


We click on the VIEW menu, 106 instances are now selected i.e. the unlabeled instances.
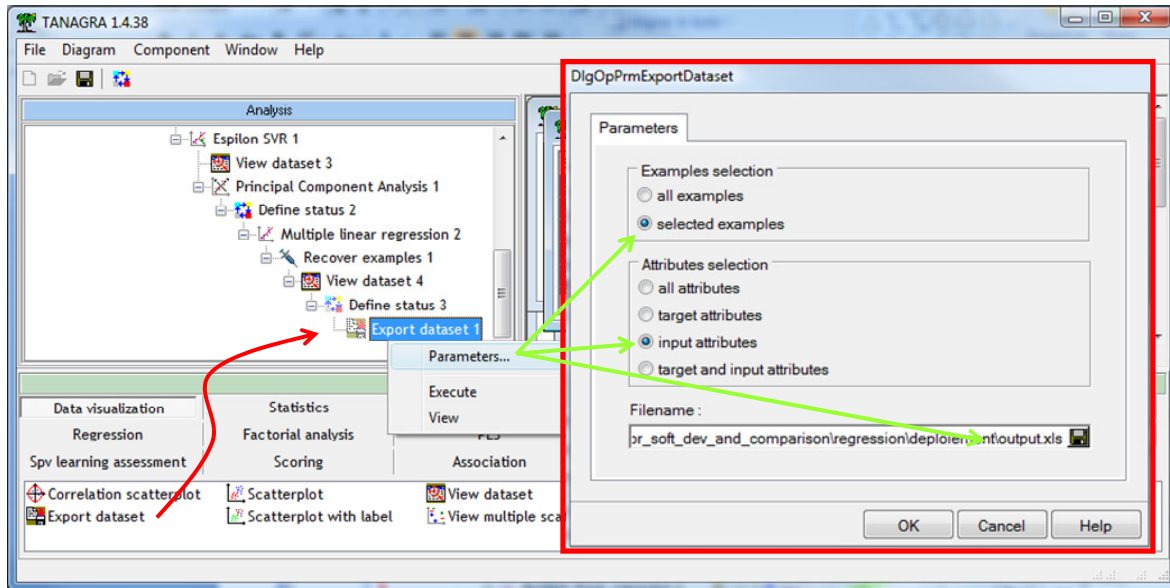
We can visualize them using the VIEW DATASET component. On the last columns, we observe the predictions provided by the various models.
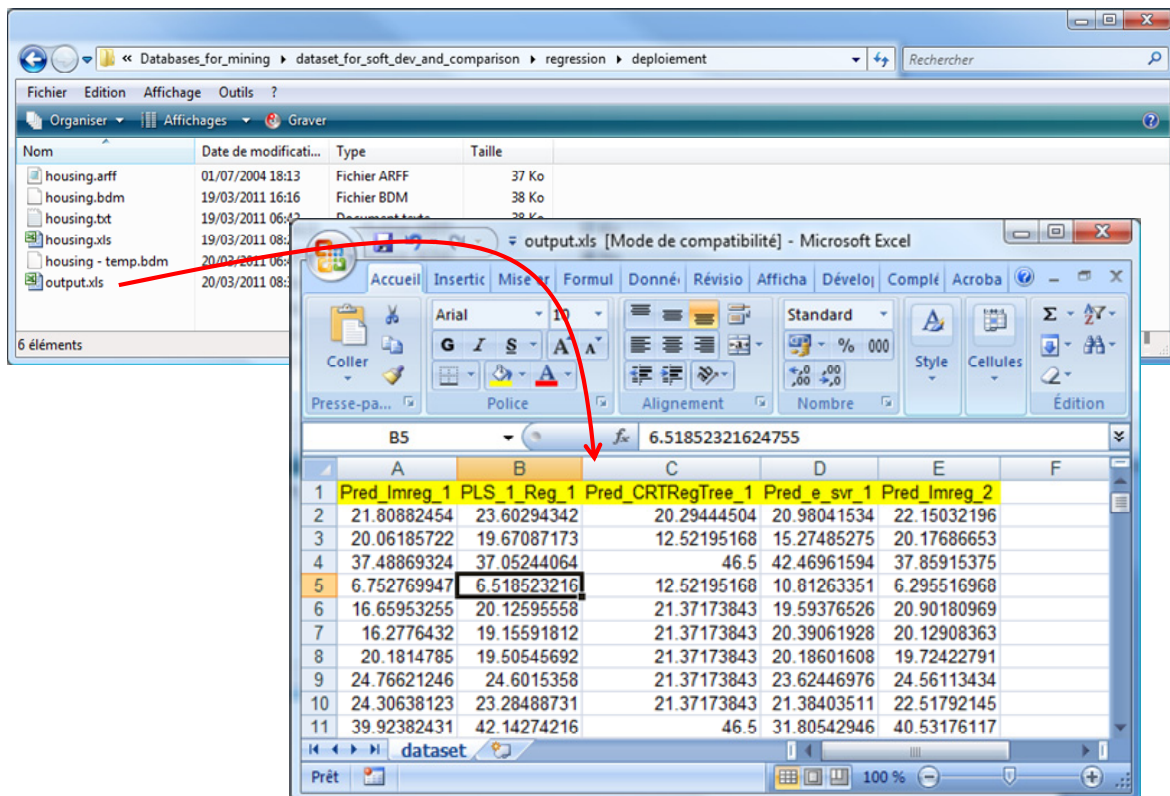


Last, we can export these values in an output file, in an Excel file format for instance. We add the DEFINE STATUS component into the diagram in order to set the predictions as INPUT.

We add the EXPORT DATASET. We export only the selected instances and the columns defined as INPUT above. The data file name is "**output.xls**".
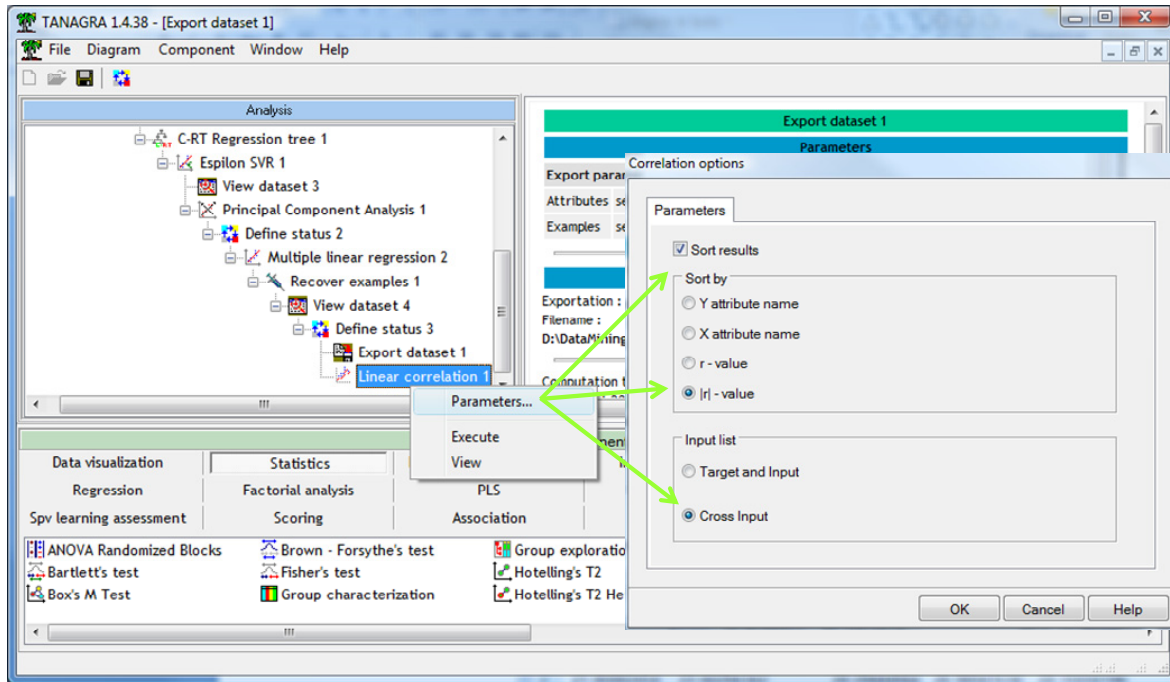


We click on the VIEW menu. The data file is rightly generated. We can use a spreadsheet application (Excel or OpenOffice/LibreOffice) to visualize them.
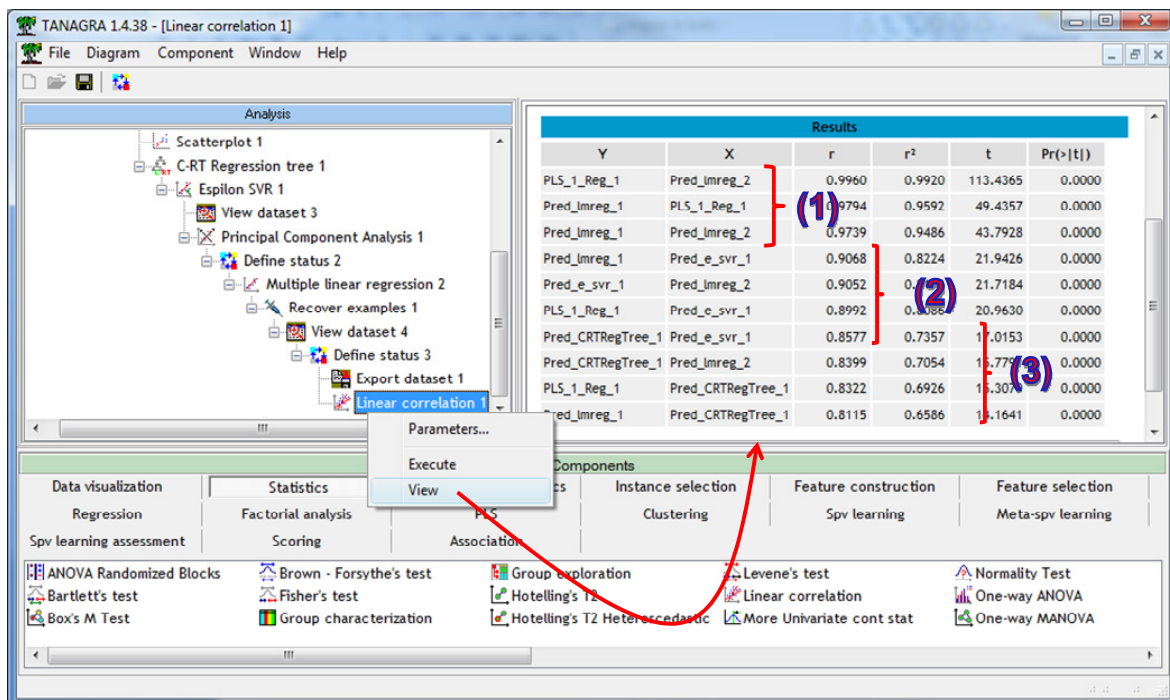


## 3.5   Comparing the predictions

Last step of our study, we want to compare the predictions of the various models. The easiest way is to calculate the correlations between them. We insert the LINEAR CORRELATION component (STATISTICS tab) into the diagram. We set the following settings. We especially want to highlight the most highly correlated predictions.

We observe that the predictions of linear models are very similar (1 - Linear regression, PLS Regression, Regression on PCA factors). The SVR with a RBF kernel, a nonlinear model, provides different predictions (2). And then, the tree which is also a nonlinear model but with different characteristics, it partitions the representation space in many regions where it makes a local prediction, it provides predictions which are not really correlated to the other ones (3).



# 4   Conclusion

Of course, other strategies exist for the regression model deployment. Especially into a professional context, when we wish to industrialize the model into a business decision-making process. A simple solution is to export the model in a standard format. The PMML format is a good example for that

(http://www.dmg.org/v4-0-1/Regression.html for the linear regression; http://www.dmg.org/v4-0-1/TreeModel.html, for the trees; http://www.dmg.org/v4-0-1/SupportVectorMachine.html for the support vector regression, we note that the model description becomes more complex; etc.). Then we use specialized tools that can handle this kind of format for the deployment (Pentaho Data Integration for instance - http://www.pentaho.com/products/data_integration/).