# Target

How to apply classifier on a new dataset?

**This functionality surpasses the TANAGRA framework, which intends only to evaluate and compare data mining algorithms. But, users ask it often; in this tutorial we show how to proceed.**

# Dataset

Data preparation is a primordial step. Indeed, TANAGRA can handle only one data source. It is not theoretically possible to manipulate two dataset, and therefore apply a classifier on a new dataset. The trick is in the dataset preparation.

We use the BREAST CANCER WINSCONCIN dataset in this tutorial (detect a tumor from cells properties). We subdivide the dataset into 500 examples for the learning phase, and 199 examples for the classification phase.

The dataset is built in several steps:

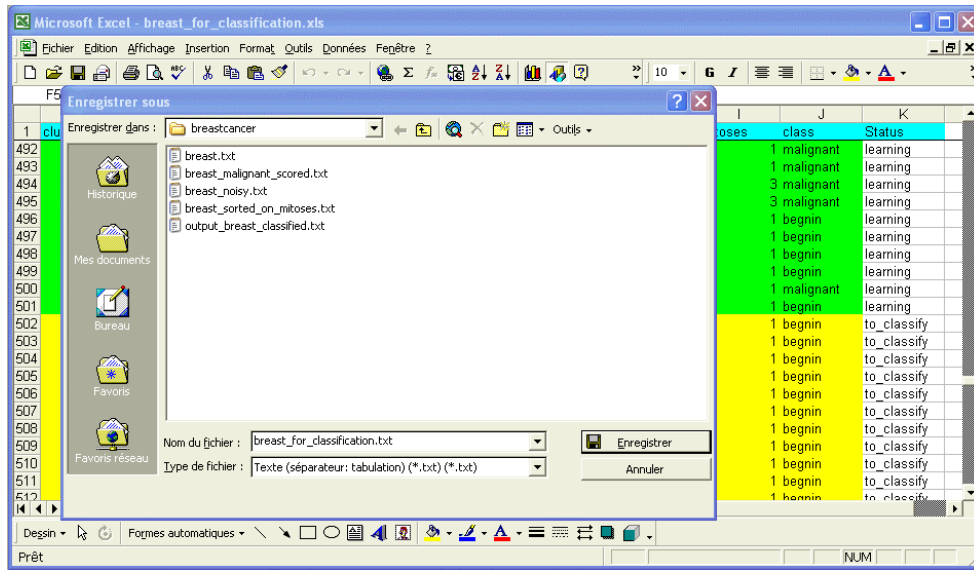(1) Join the two dataset in a file, using a spreadsheet for instance, you can set the learning set before the examples to classify.



(2) Add a new attribute "STATUS" with two values "learning" and "To_classifiy" which allows us to specify the status of each example.

(3) At last, even if it is not very intuitive, you must assign a class value for each example to classify. The reason is that TANAGRA does not handle missing data. You must use one of the existing values. This information will not be used in the following. In our dataset, we set the 199 examples to classify to "begnin".
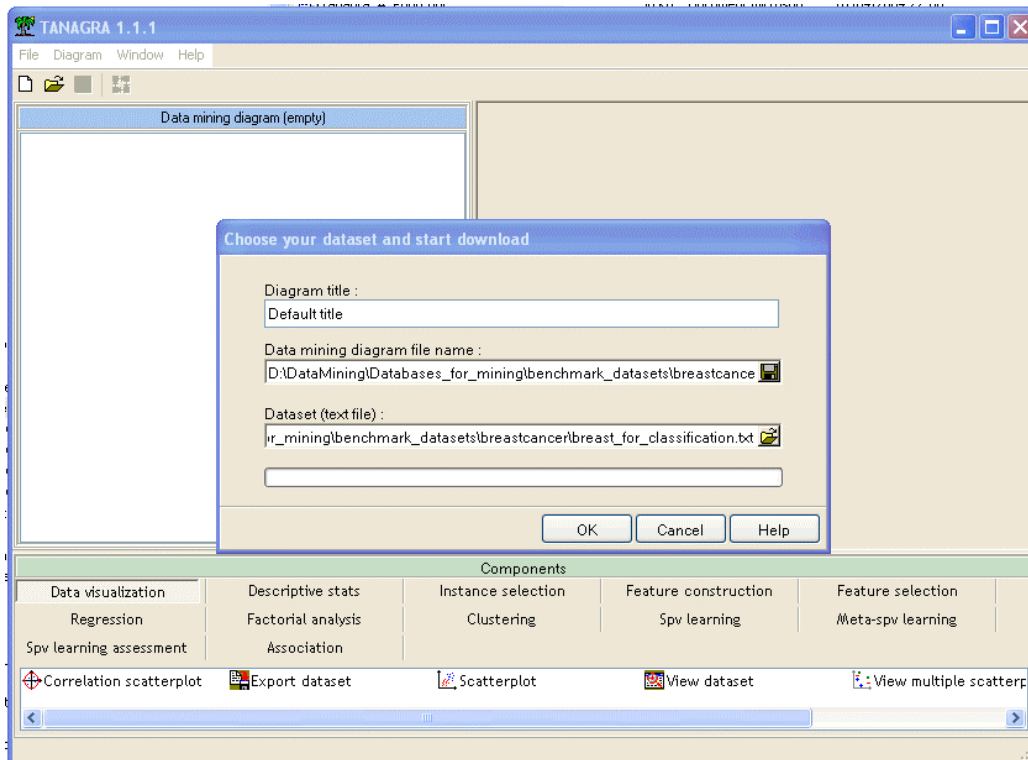
(4) Export the dataset in a tab separator text file for TANAGRA.
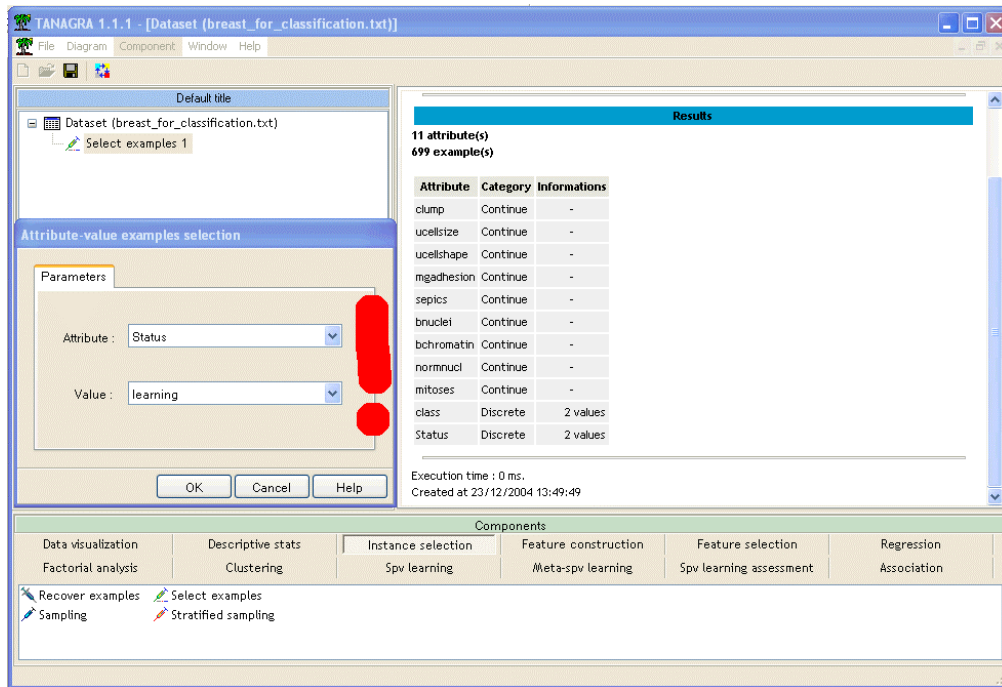


# Classifier deployment on a new dataset

## Dataset importation in TANAGRA

Import the dataset in TANAGRA and define a stream diagram.

## Select the learning set

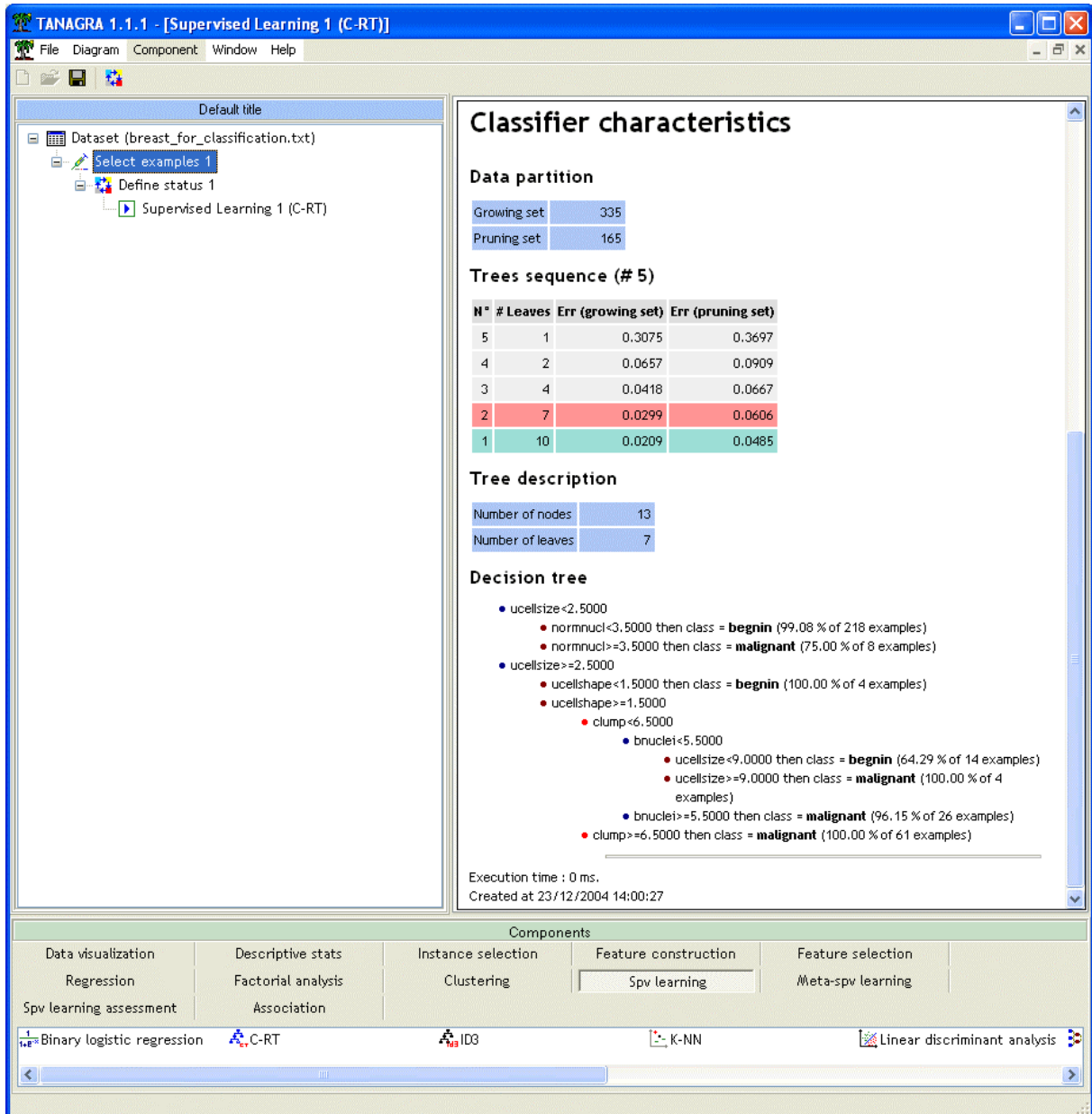Above all, you must select the learning set. Use SELECT EXAMPLES component using the STATUS attribute: we have 500 examples for the learning phase (active examples), and 199 examples for the classification phase (idle examples).



## Supervised learning

In the next step, we must select examples and define the learning algorithm. We use the Breiman's et al. (1984) famous classification tree algorithm.
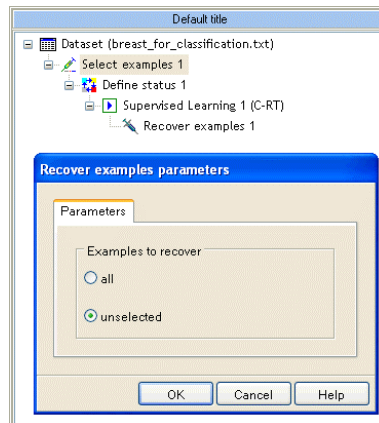
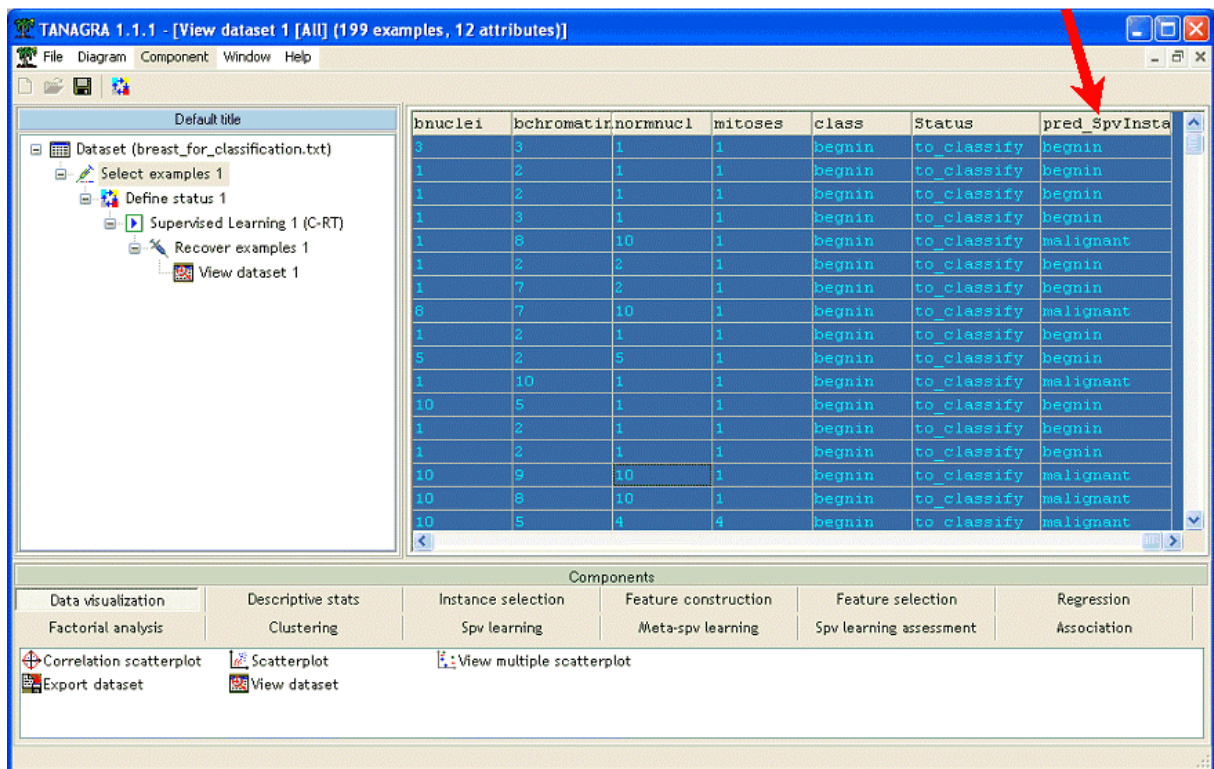TANAGRA adds automatically a new attribute, the predicted values.

**Note that if the learning phase is realized on the learning set, the classification is realized on the whole examples, including the idle examples. We exploit this property to classify new dataset or to test the classifier on an external test sample.**

## View classified examples

We can view the predicted attribute on the examples to classify using a VIEW DATASET component. First, we must set idle examples to active examples with RECOVER EXAMPLES component.
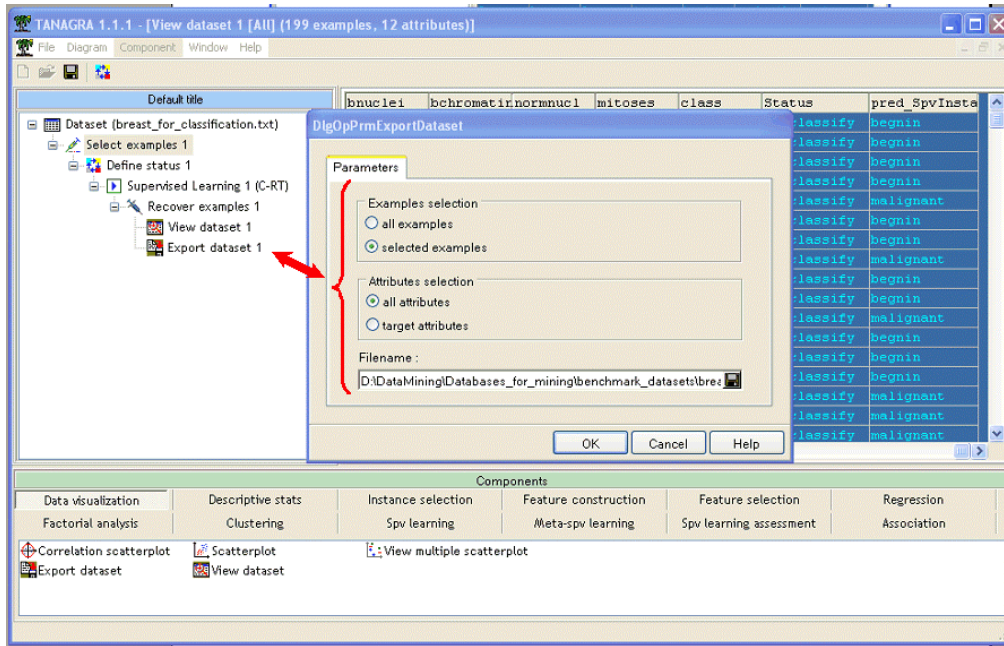
VIEW DATASET



## Export results

The last step is to export the examples with EXPORT DATASET component. We choose to export all attributes on examples to classify.

We can view the dataset in a spreadsheet.

# Conclusion

We can follow the same way for classifier validation on an external test set. To obtain classification matrix, you can build a contingency table between class attribute and predicted attribute on the idle examples.

I know this method is very complex but it is the only way to apply a classifier on a new dataset.

For classifier validation with subset of examples which not used for the learning phase, you can also use SAMPLING or ASSESSMENT components.