

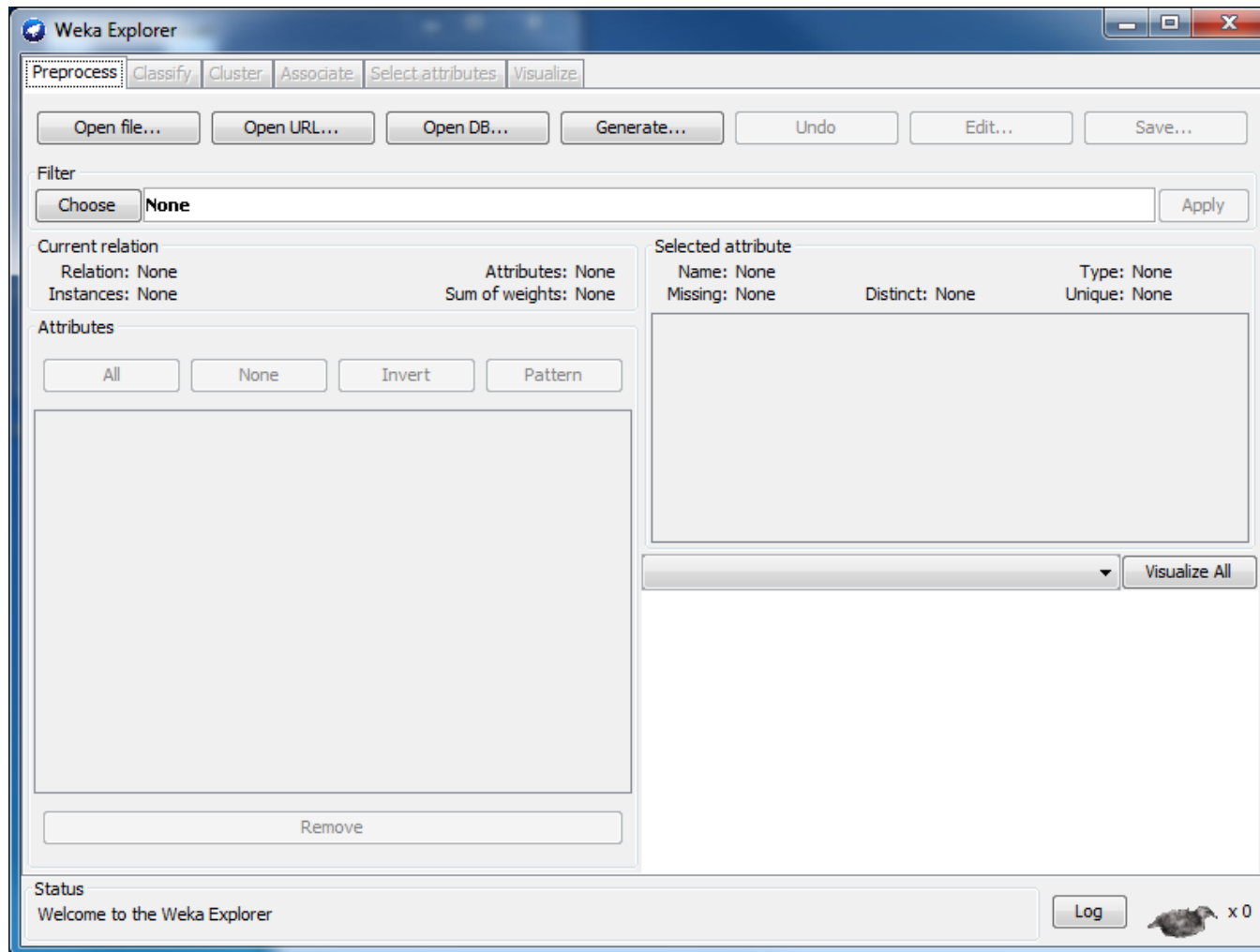


Data mining with WEKA

A use-case to help you get started

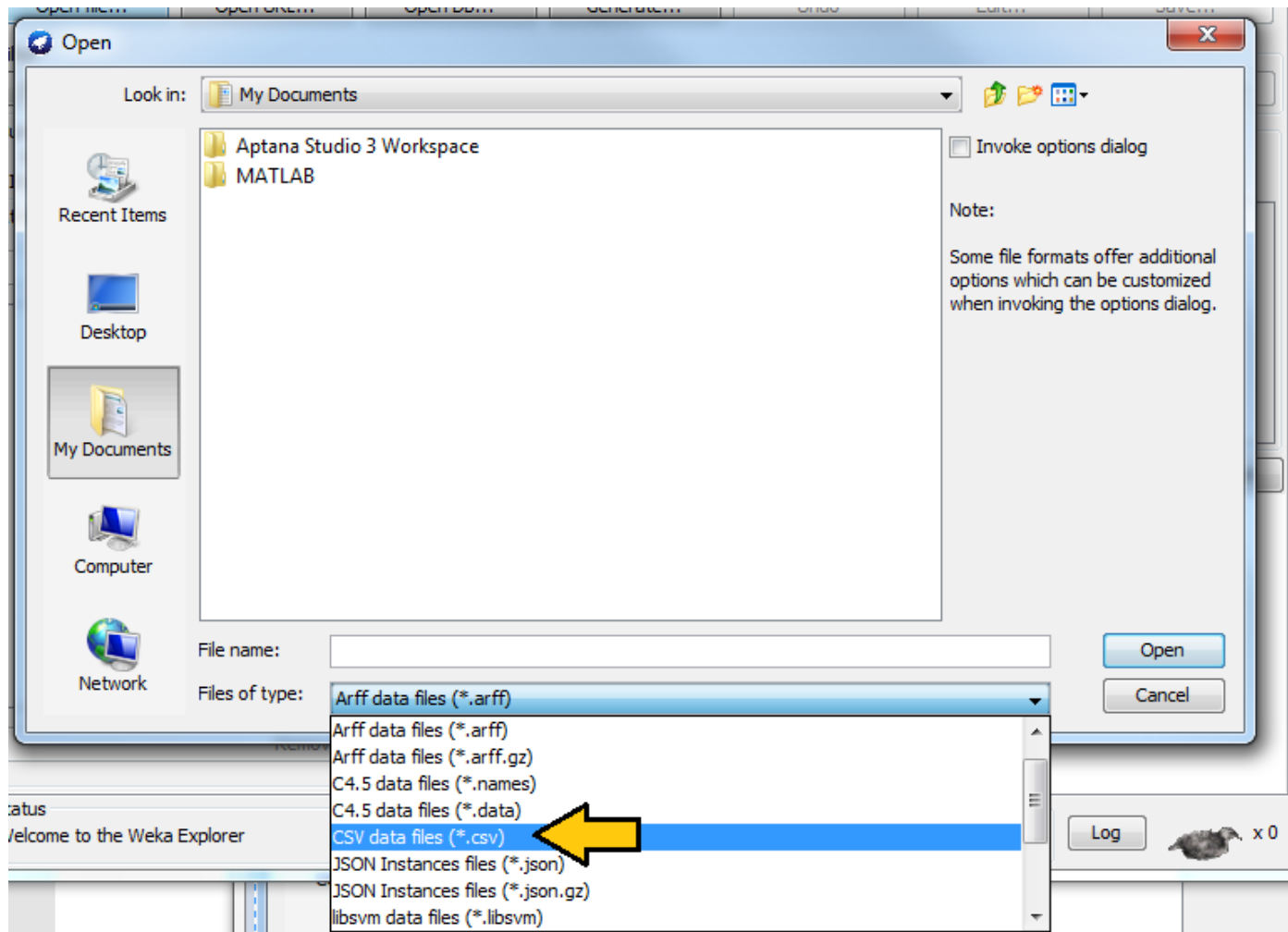
Charalampos Mavroforakis
BU CS105, Fall 2011

Starting WEKA



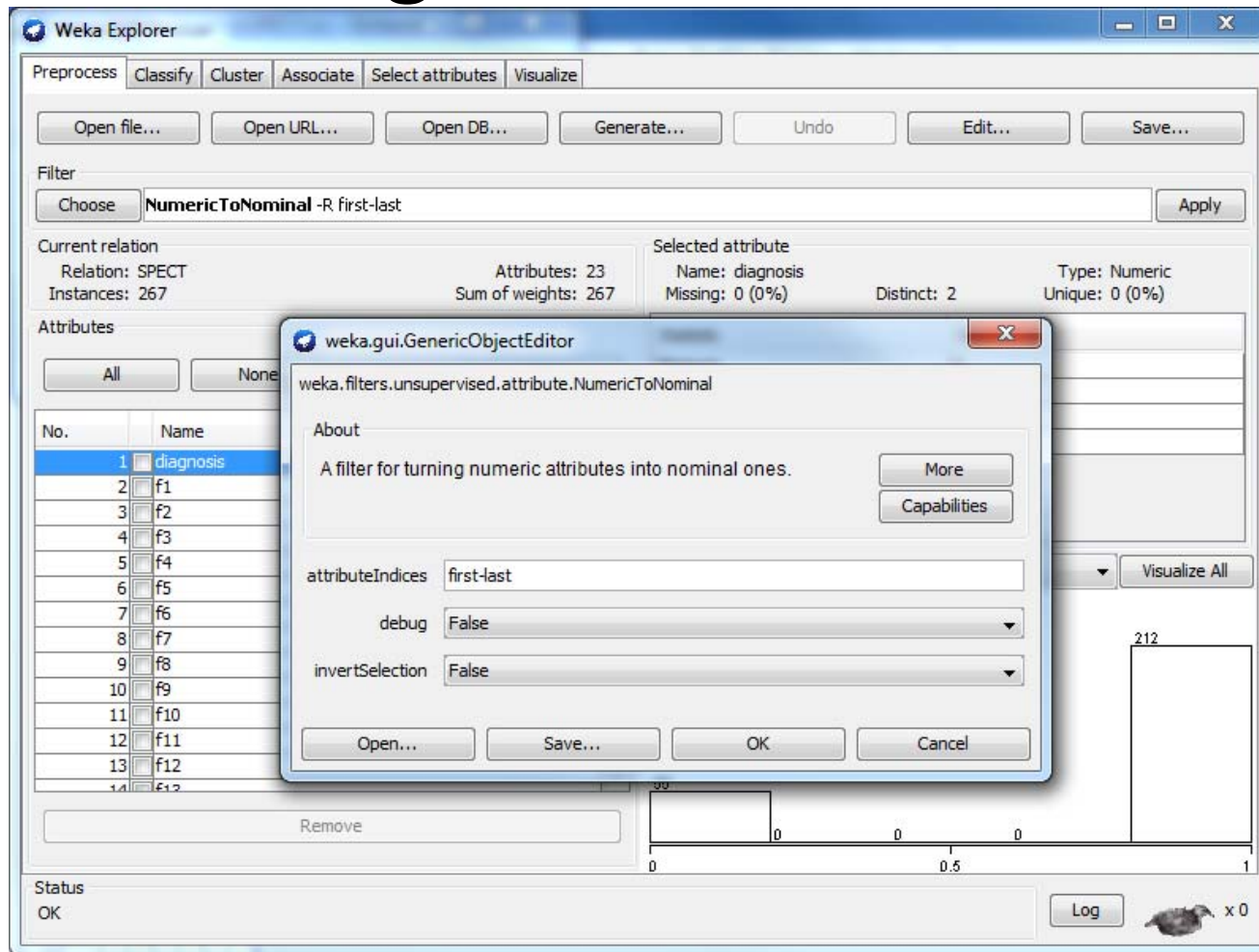
Open Weka : Start > All Programs > Weka 3.x.x > Weka 3.x
From the "*Weka GUI Chooser*", pick "*Explorer*". This is the main WEKA tool that we are going to use.

Opening a dataset



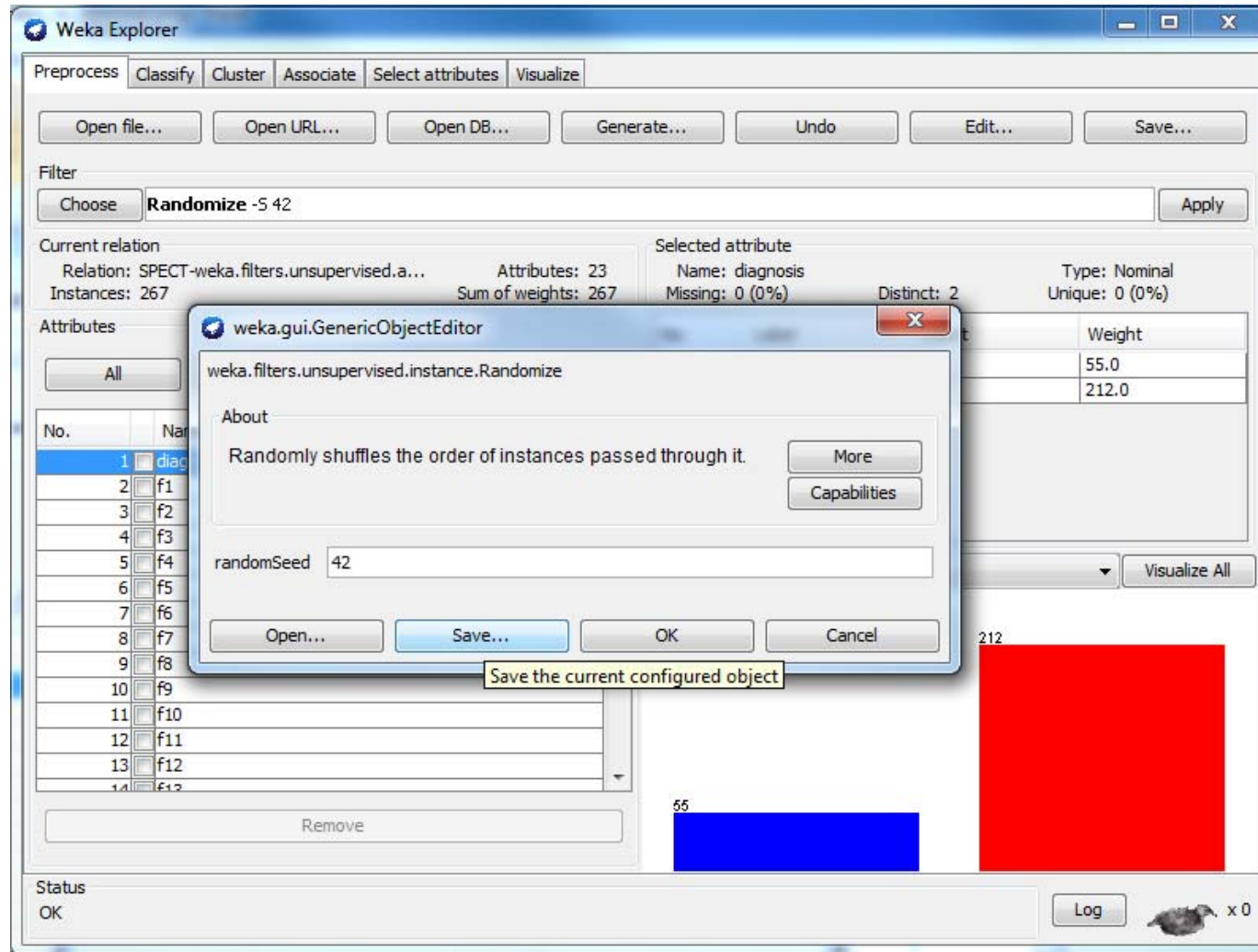
To open a dataset (a .csv file in our case), we click "Open file ..." in the *Preprocess* tab and open the file that contains our data. **Remember** that in the open menu you have to choose csv if your file was saved as such. Let's open SPECT.csv

Transforming values to nominal (if needed)



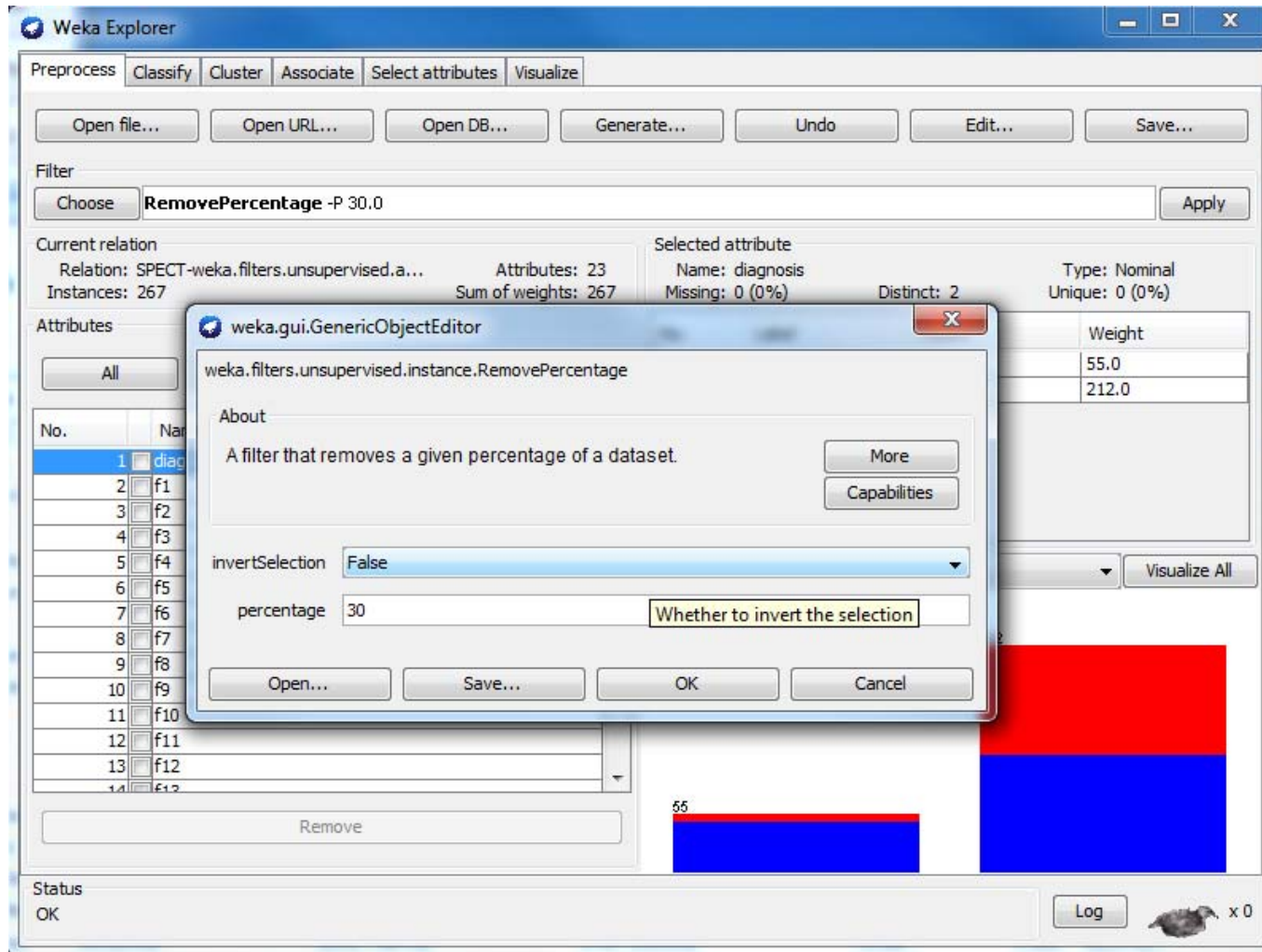
Weka classified every attribute in our dataset as numeric, so we have to manually transform them to nominal. To do so, we will use a filter. We navigate to *NumericToNominal*, which is in *Unsupervised > attribute*. If we click on that, we will get to the options of that filter. Mainly, the most interesting one here is the *attributeIndices*, which enumerates all the attributes that you want the filter to be applied on. To finish, we click *Apply*.

Splitting the dataset



We have to split the dataset into two, 30% testing and 70% training. To do that, we first *Randomize* the dataset (Unsupervised > Instance), so that we create a random permutation.

Splitting the dataset



Then we apply *RemovePercentage* (Unsupervised > Instance) with percentage 30 and save the resulting dataset as training.

Splitting the dataset

The screenshot shows the Weka Explorer interface with the 'Preprocess' tab selected. A 'RemovePercentage' filter is applied with a percentage of 30.0. A dialog box titled 'weka.gui.GenericObjectEditor' is open, showing the configuration for the filter. The 'invertSelection' dropdown is set to 'True', indicated by a yellow arrow. The 'percentage' field is set to 30.0. The dialog also includes an 'About' section and buttons for 'Open...', 'Save...', 'OK', and 'Cancel'.

Current relation: SPECT-weka.filters.unsupervised.a...
Attributes: 23
Instances: 187
Sum of weights: 187
Selected attribute: Name: diagnosis
Missing: 0 (0%)
Distinct: 2
Type: Nominal
Unique: 0 (0%)

Weight
32.0
155.0

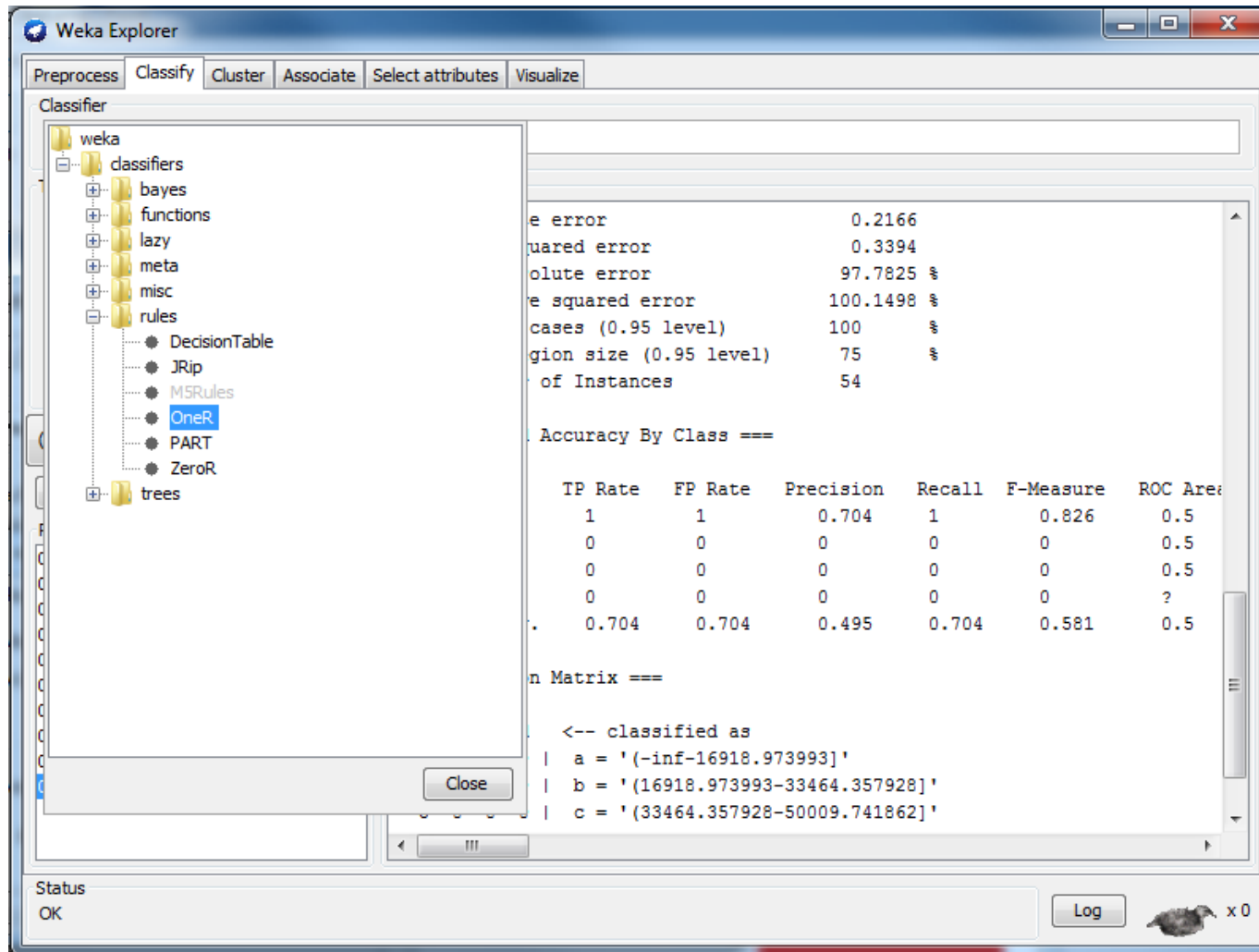
invertSelection: True
percentage: 30.0

Remove

Status: OK

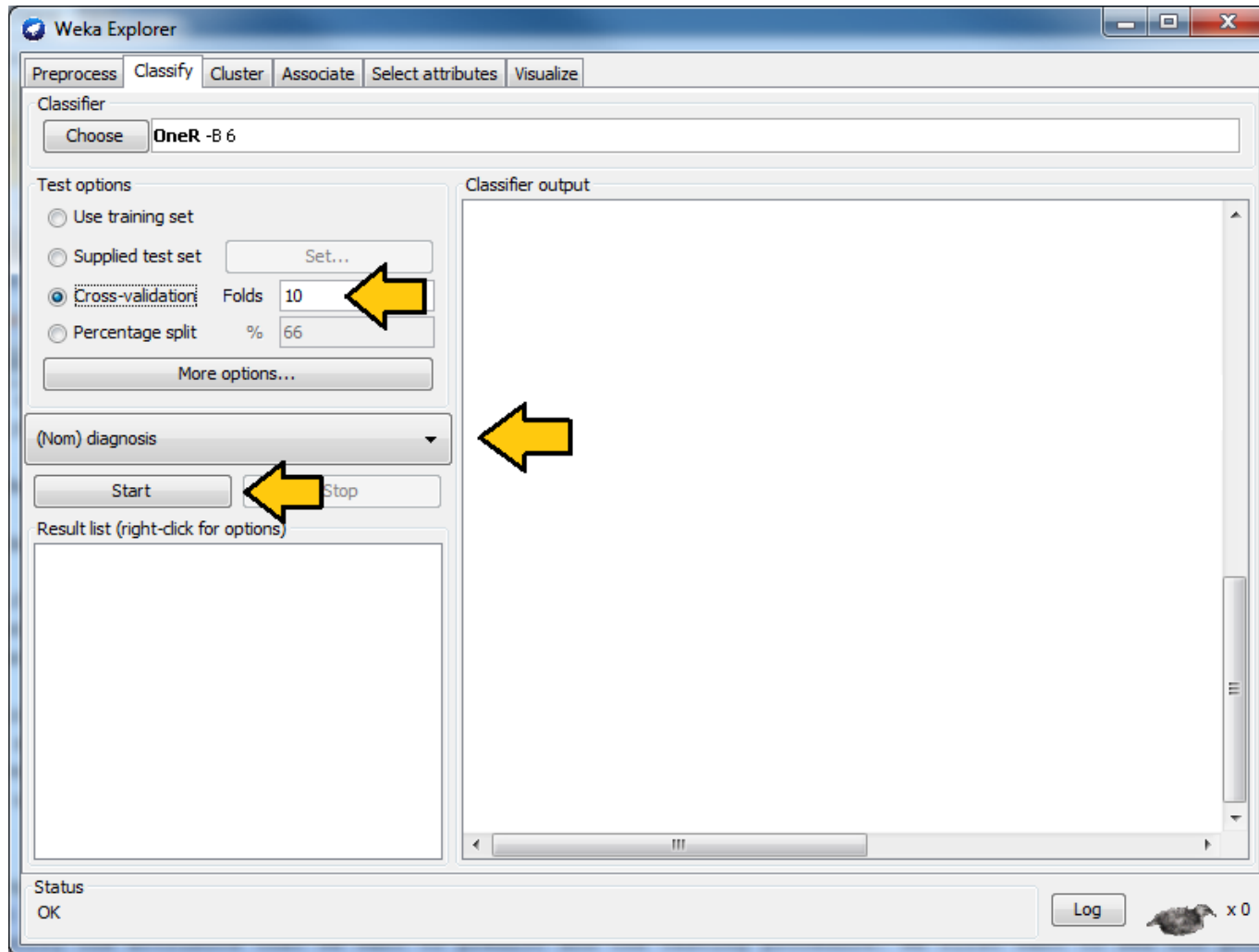
After that, we undo and apply the same filter choosing *invertSelection* this time. This will pick the rest of the data (30%) so we save them as the testing.

Training models



From now on we will be using the training dataset. We switch to the tab "Classify" and we pick a classifier. Let's start with *OneR*, which is the same with the one we saw in the class.

Training models



We have to specify the attribute that we want to predict and the testing procedure. We first want to see how good OneR is as a model, so we use cross-validation. , and only after that will we go and check what it predicts on the unseen data.

Training models

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose OneR -B 6

Test options

Use training set

Supplied test set Set...

Cross-validation **Test on a user-specified dataset**

Percentage split % 66

More options...

(Nom) diagnosis

Start Stop

Result list (right-click for options)

02:00:07 - rules.OneR

Classifier output

```
Correctly Classified Instances      149      79.6791 %
Incorrectly Classified Instances    38      20.3209 %
Kappa statistic                      0
Absolute error                      0.2032
Root mean squared error              0.4508
Relative absolute error              62.3438 %
Root relative squared error          111.9988 %
Coverage of cases (0.95 level)      79.6791 %
Mean rel. region size (0.95 level)  50 %
Total Number of Instances           187

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  ROC Area
          0         0         0           0         0           0.5
          1         1         0.797       1         0.887       0.5
Weighted Avg.   0.797   0.797   0.635       0.797   0.707       0.5

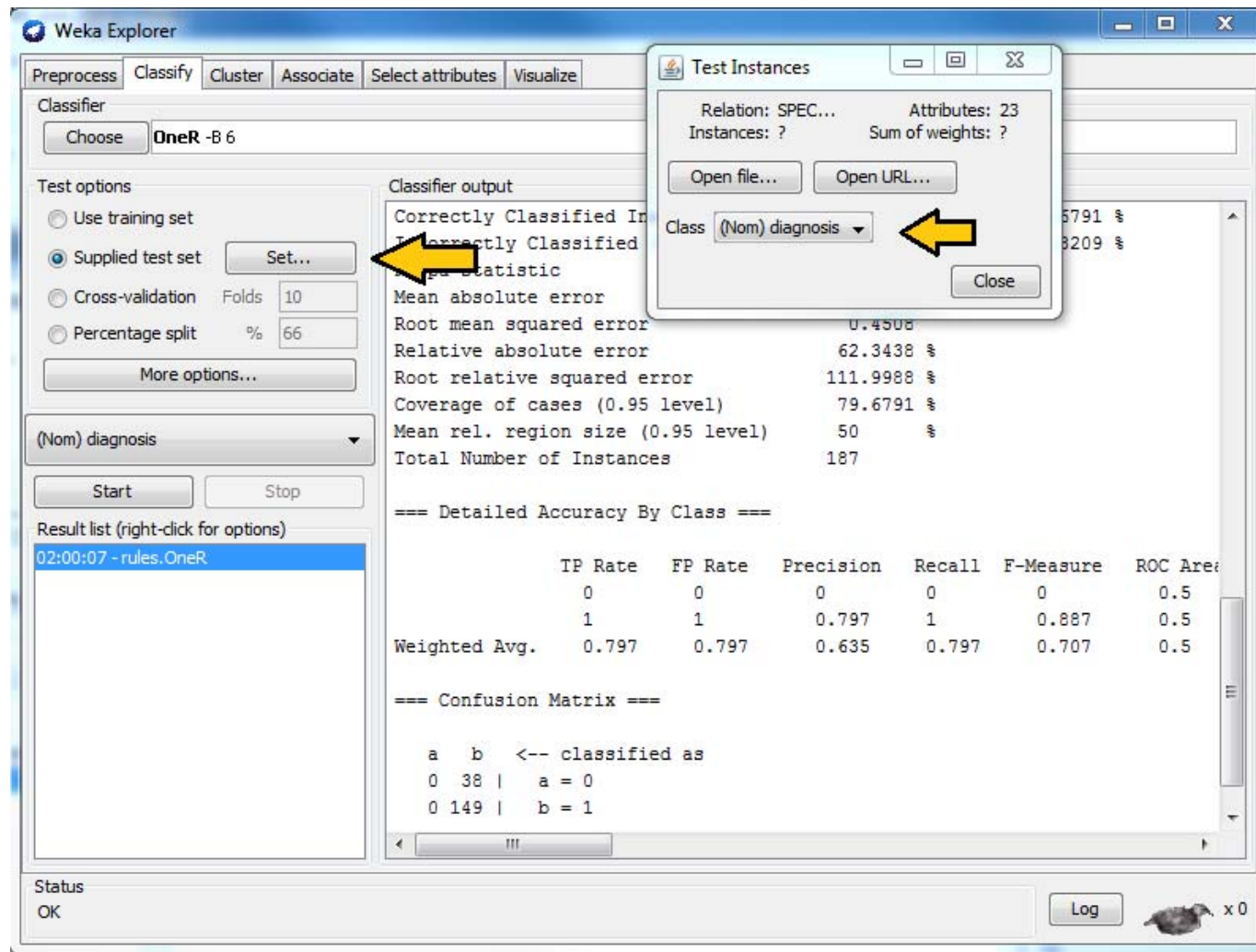
=== Confusion Matrix ===

 a  b  <-- classified as
 0  38 |  a = 0
 0 149 |  b = 1
```

Status OK Log x 0

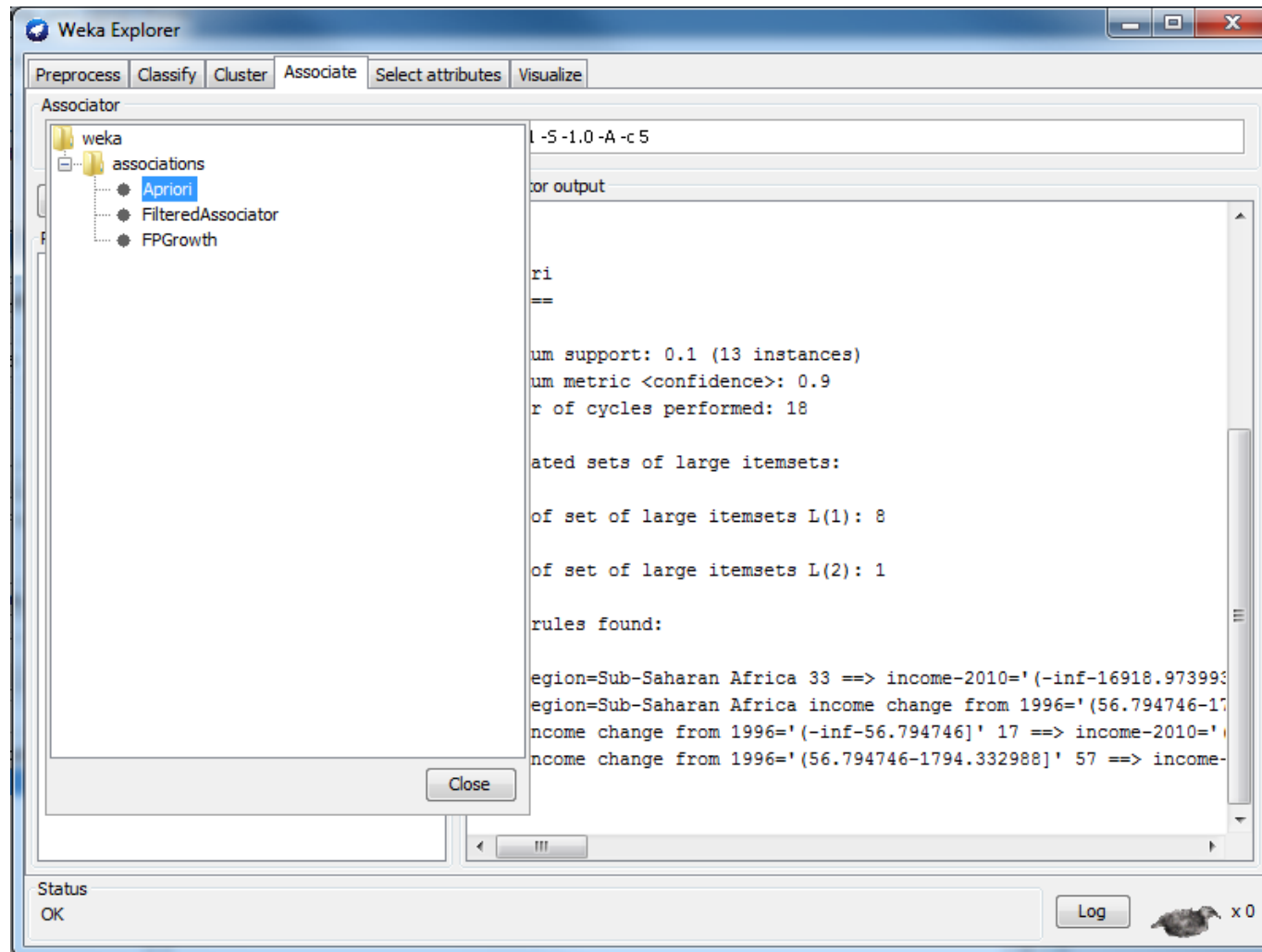
In the output, we get information about the average accuracy and the confusion matrix of our model.

Training models



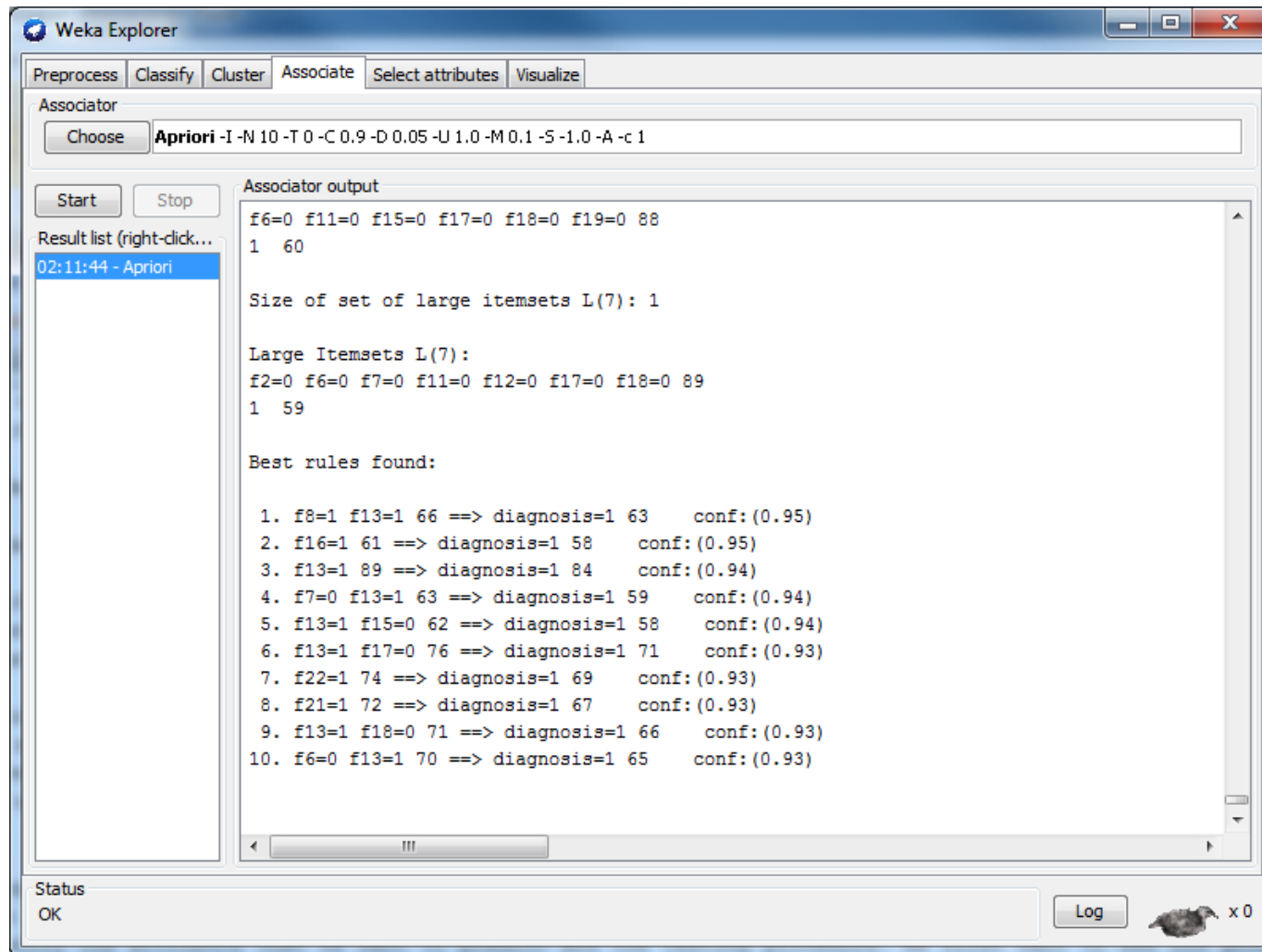
In order to check how well we do on the unseen data, we select "*supplied test set*", we open the testing dataset that we have created and we specify which attribute is the class. We run the algorithm again and we notice the differences in the confusion matrix and the accuracy.

Association learning



If all of our attributes are nominal (in case they are not, we can discretize them in the *Preprocess* tab) we can also do association learning. In order to do that, we switch to the *Association* tab and we choose the *Apriori* algorithm. You can play around with its parameters if you want.

Association learning



The screenshot shows the Weka Explorer interface with the 'Associate' tab selected. The 'Associator' dropdown is set to 'Apriori'. The command line shows: `Apriori -I -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -A -c 1`. The 'Associator output' pane displays the following text:

```
f6=0 f11=0 f15=0 f17=0 f18=0 f19=0 88
1 60

Size of set of large itemsets L(7): 1

Large Itemsets L(7):
f2=0 f6=0 f7=0 f11=0 f12=0 f17=0 f18=0 89
1 59

Best rules found:

1. f8=1 f13=1 66 ==> diagnosis=1 63   conf:(0.95)
2. f16=1 61 ==> diagnosis=1 58   conf:(0.95)
3. f13=1 89 ==> diagnosis=1 84   conf:(0.94)
4. f7=0 f13=1 63 ==> diagnosis=1 59   conf:(0.94)
5. f13=1 f15=0 62 ==> diagnosis=1 58   conf:(0.94)
6. f13=1 f17=0 76 ==> diagnosis=1 71   conf:(0.93)
7. f22=1 74 ==> diagnosis=1 69   conf:(0.93)
8. f21=1 72 ==> diagnosis=1 67   conf:(0.93)
9. f13=1 f18=0 71 ==> diagnosis=1 66   conf:(0.93)
10. f6=0 f13=1 70 ==> diagnosis=1 65   conf:(0.93)
```

The 'Result list' on the left shows a single entry: '02:11:44 - Apriori'. The 'Status' bar at the bottom indicates 'OK'.

We could set *car* to True (so that it produces rules that predict the class attribute) and specify the index of the attribute that will be considered as class. *minMetric* sets the threshold of confidence and *numRules* limits the number of rules that will be created. The result will be a set of rules that predict the class, together with their confidence.