

DMAT – Workbook 03

Classification via Decision Trees in WEKA

The following guide is based WEKA version 3.4.1 Additional resources on WEKA, including sample data sets can be found from the official [WEKA Web site](http://www.cs.waikato.ac.nz/ml/weka/).

This example illustrates the use of C4.5 (J48) classifier in WEKA. The sample data set used for this example, unless otherwise indicated, is the bank data available in comma-separated format ([bank-data.csv](#)). This document assumes that appropriate data preprocessing has been performed. In this case ID field has been removed. Since C4.5 algorithm can handle numeric attributes, there is no need to discretize any of the attributes. For the purposes of this example, however, the "Children" attribute has been converted into a categorical attribute with values "YES" or "NO".

WEKA has implementations of numerous classification and prediction algorithms. The basic ideas behind using all of these are similar. In this example we will use the **modified** version of the bank data to classify new instances using the C4.5 algorithm (note that the C4.5 is implemented in WEKA by the classifier class: `weka.classifiers.trees.J48`). The modified (and smaller) version of the bank data can be found in the file “bank.arff” and the new unclassified instances are in the file “bank-new.arff”.

NB The datasets provided differ slightly from those in the diagrams.

- bank.arff – 500 instances
- bank-new.arff has 100 instances

As usual, we begin by loading the data (“bank.arff”) into WEKA, as seen in Figure 20:

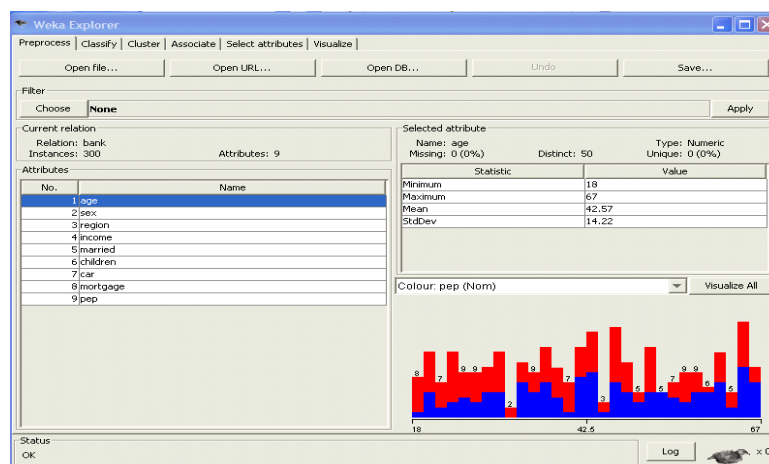
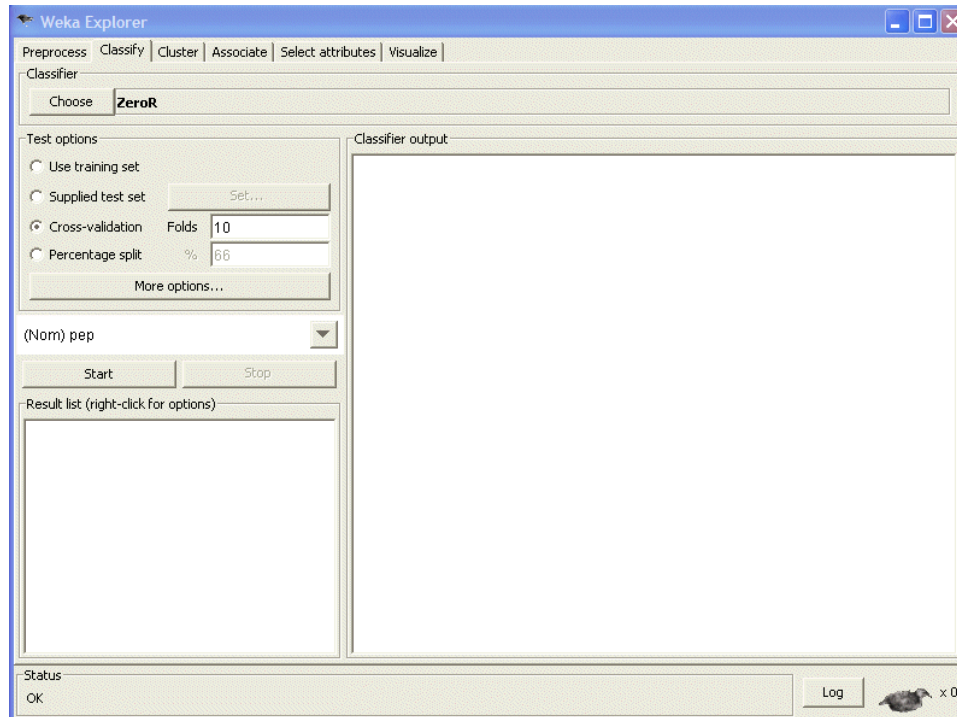
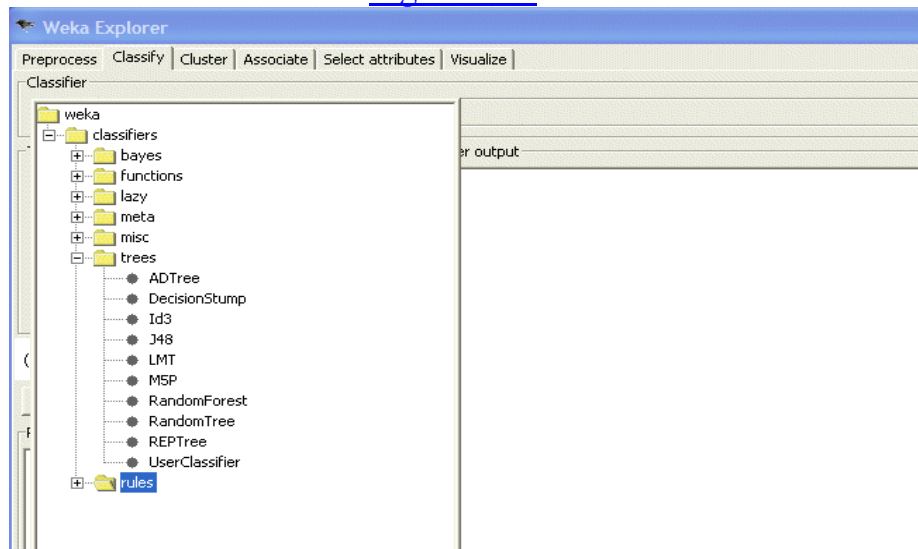


Figure 20

Next, we select the "Classify" tab and click the "Choose" button to select the J48 classifier, as depicted in Figures 21-a and 21-b. Note that J48 (implementation of C4.5 algorithm) does not require discretization of numeric attributes, in contrast to the ID3 algorithm from which C4.5 has evolved.

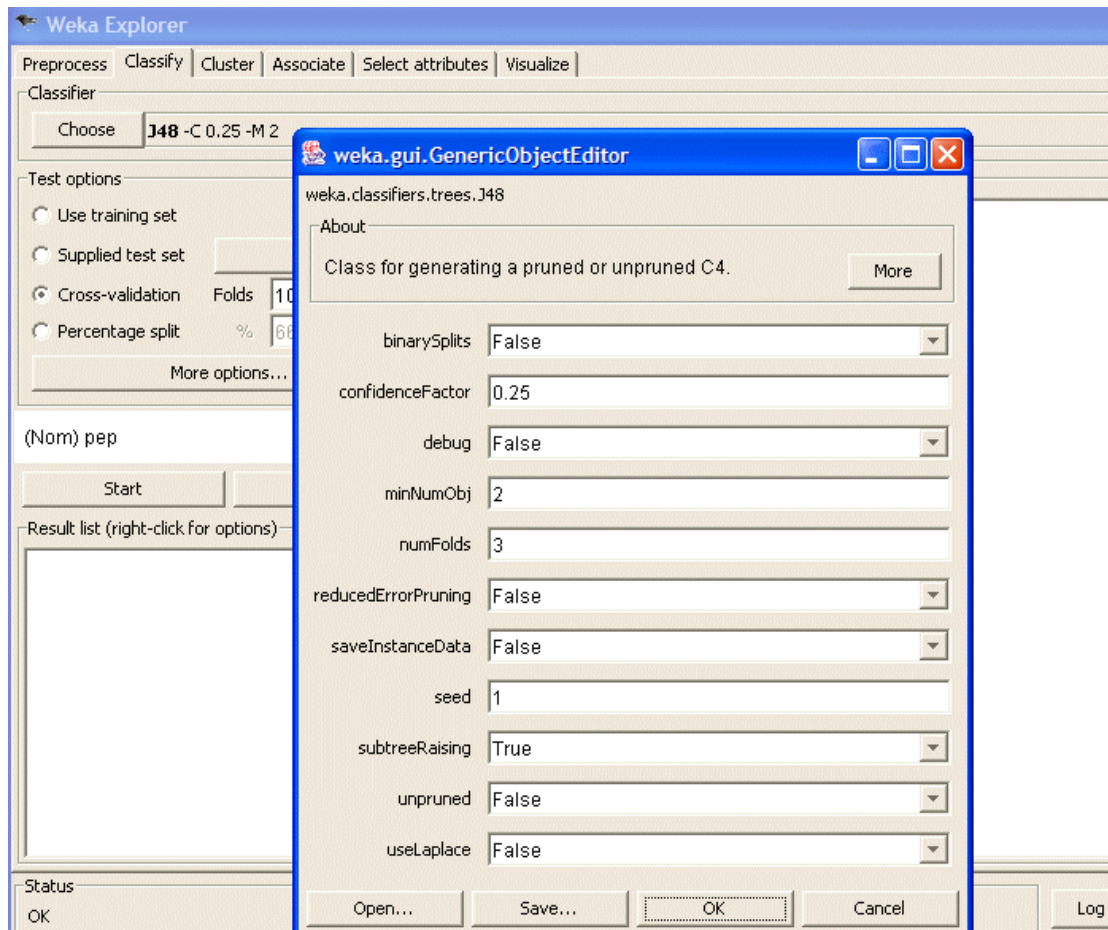


[Figure 21-a](#)



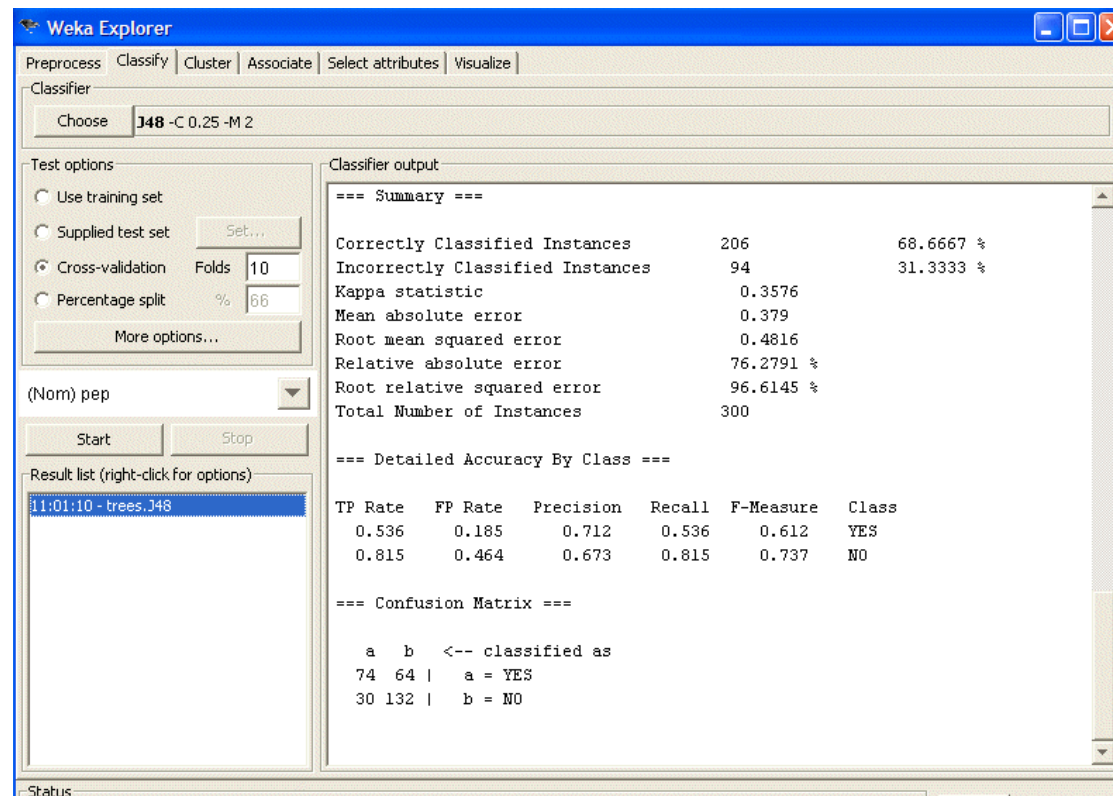
[Figure 21-b](#)

Now, we can specify the various parameters. These can be specified by clicking in the text box to the right of the "Choose" button, as depicted in Figure 22. In this example we accept the default values. The default version does perform some pruning (using the subtree raising approach), but does not perform error pruning. The selected parameters are depicted in Figure 22.



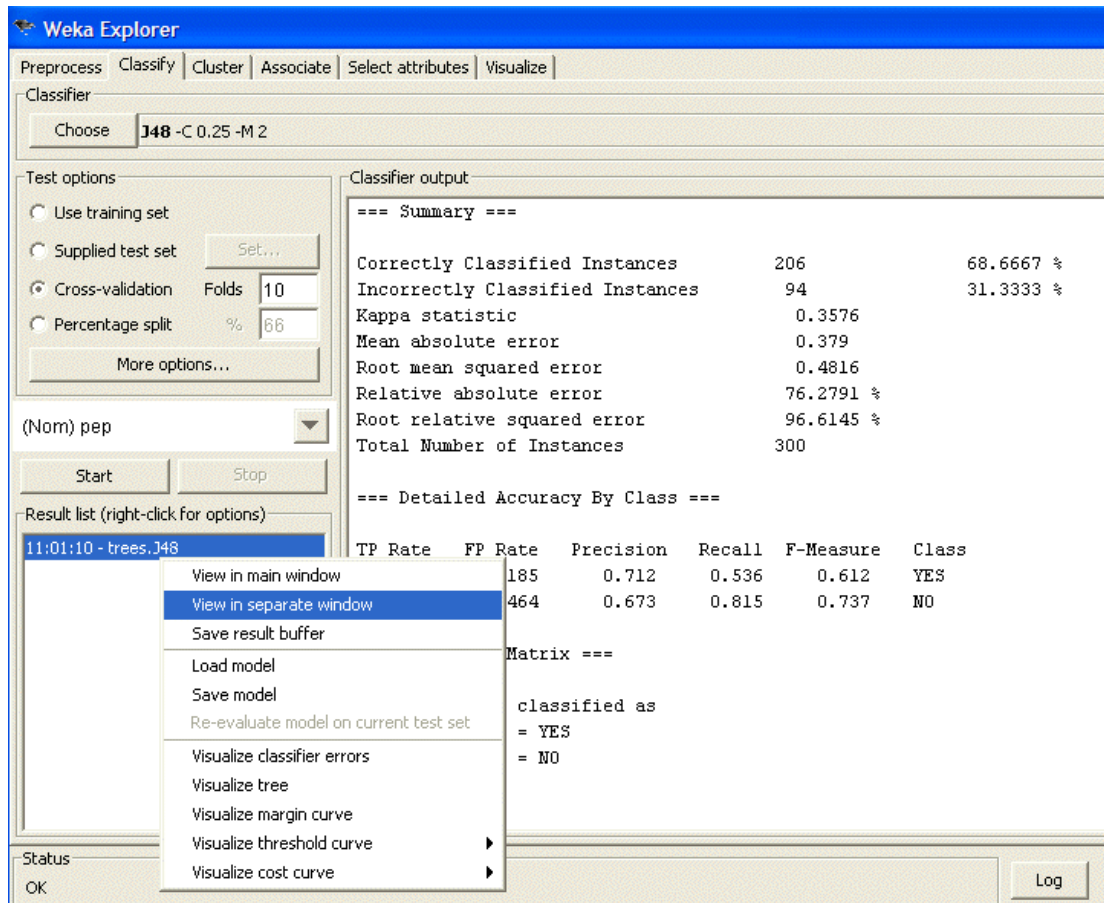
[Figure 22](#)

Under the "Test options" in the main panel we select 10-fold cross-validation as our evaluation approach. Since we do not have separate evaluation data set, this is necessary to get a reasonable idea of accuracy of the generated model. We now click "Start" to generate the model. The ASCII version of the tree as well as evaluation statistics will appear in the eight panel when the model construction is completed (see Figure 23).



[Figure 23](#)

We can view this information in a separate window by right clicking the last result set (inside the "Result list" panel on the left) and selecting "View in separate window" from the pop-up menu. These steps and the resulting window containing the classification results are depicted in Figures 24-a and 24-b.



[Figure 24-a](#)

```
11:01:10 - trees.J48

Test mode:      10-fold cross-validation

=== Classifier model (full training set) ===

J48 pruned tree
-----

children = YES
|  income <= 30099.3
|  |  car = YES: NO (50.0/15.0)
|  |  car = NO
|  |  |  married = YES
|  |  |  |  income <= 13106.6: NO (9.0/2.0)
|  |  |  |  income > 13106.6
|  |  |  |  |  mortgage = YES: YES (12.0/3.0)
|  |  |  |  |  mortgage = NO
|  |  |  |  |  income <= 18923: YES (9.0/3.0)
|  |  |  |  |  income > 18923: NO (10.0/3.0)
|  |  |  |  |  married = NO: NO (22.0/6.0)
|  |  |  |  income > 30099.3: YES (59.0/7.0)
children = NO
|  married = YES
|  |  mortgage = YES
|  |  |  region = INNER_CITY
|  |  |  |  income <= 39547.8: YES (12.0/3.0)
|  |  |  |  income > 39547.8: NO (4.0)
|  |  |  region = RURAL: NO (3.0/1.0)
|  |  |  region = TOWN: NO (9.0/2.0)
|  |  |  region = SUBURBAN: NO (4.0/1.0)
|  |  mortgage = NO: NO (57.0/9.0)
|  married = NO
|  |  mortgage = YES
|  |  |  age <= 39
|  |  |  |  age <= 28: NO (4.0)
|  |  |  |  age > 28: YES (5.0/1.0)
|  |  |  age > 39: NO (11.0)
|  |  mortgage = NO: YES (20.0/1.0)

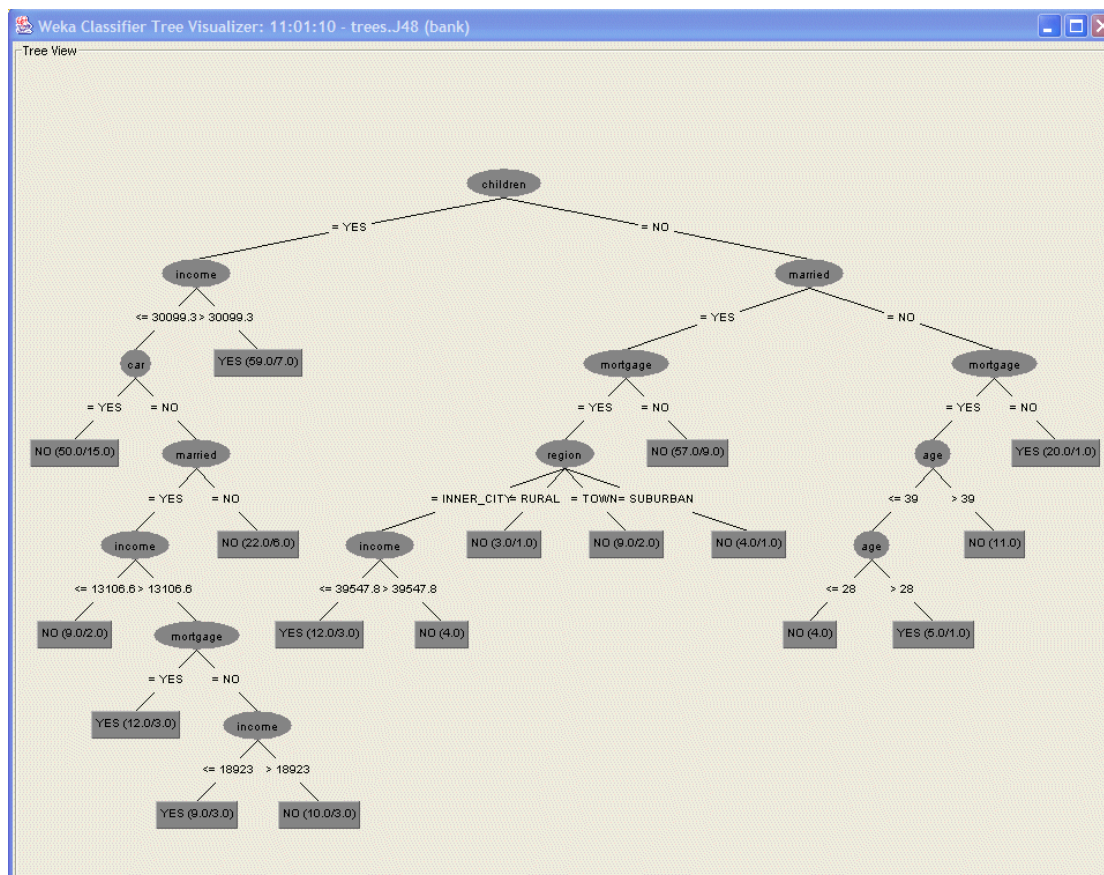
Number of Leaves :      17

Size of the tree :      31
```

[Figure 24-b](#)

Note that the classification accuracy of our model is only about 88%. This may indicate that we may need to do more work (either in preprocessing or in selecting the correct parameters for classification), before building another model. In this example, however, we will continue with this model despite its inaccuracy.

WEKA also let's us view a graphical rendition of the classification tree. This can be done by right clicking the last result set (as before) and selecting "Visualize tree" from the pop-up menu. The tree for this example is depicted in Figure 25. Note that by resizing the window and selecting various menu items from inside the tree view (using the right mouse button), we can adjust the tree view to make it more readable.



[Figure 25](#)

We will now use our model to classify the new instances. A portion of the new instances ARFF file is depicted in Figure 26. Note that the attribute section is identical to the training data (bank data we used for building our model). However, in the data section, the value of the "pep" attribute is "?" (or unknown).

```

@relation bank-new.csv

@attribute age numeric
@attribute sex {MALE,FEMALE}
@attribute region {INNER_CITY,RURAL,TOWN,SUBURBAN}
@attribute income numeric
@attribute married {YES,NO}
@attribute children {YES,NO}
@attribute car {YES,NO}
@attribute mortgage {YES,NO}
@attribute pep {YES,NO}

@data
23,MALE,INNER_CITY,18766.9,YES,NO,YES,YES,?
30,MALE,RURAL,9915.67,NO,YES,NO,YES,?
45,FEMALE,RURAL,21881.6,NO,NO,YES,NO,?
50,MALE,TOWN,46794.4,YES,YES,NO,YES,?
41,FEMALE,INNER_CITY,20721.1,YES,NO,YES,NO,?
20,MALE,INNER_CITY,16688.5,NO,YES,NO,YES,?
46,FEMALE,RURAL,39068,YES,NO,YES,YES,?
50,FEMALE,INNER_CITY,27740.8,YES,YES,YES,YES,?
42,MALE,INNER_CITY,33584.9,NO,YES,YES,NO,?
57,FEMALE,TOWN,19621.3,YES,YES,YES,NO,?
63,FEMALE,INNER_CITY,47630.9,YES,NO,NO,YES,?
26,FEMALE,INNER_CITY,22378.5,NO,NO,YES,YES,?
62,FEMALE,RURAL,20837.1,YES,NO,YES,NO,?
26,FEMALE,SUBURBAN,23912.7,YES,NO,YES,NO,?

```

Figure 26

In the main panel, under "Test options" click the "Supplied test set" radio button, and then click the "Set..." button. This will pop up a window which allows you to open the file containing test instances, as in Figures 27-a and 27-b.

Test Instances

Relation: None
Instances: None
Attributes: None

Open file... Open URL...

Correctly Classified Instances: 206 68.6667 %
Incorrectly Classified Instances: 94 31.3333 %
Kappa statistic: 0.3576
Mean absolute error: 0.379
Root mean squared error: 0.4816
Relative absolute error: 76.2791 %
Root relative squared error: 96.6145 %
Total Number of Instances: 300

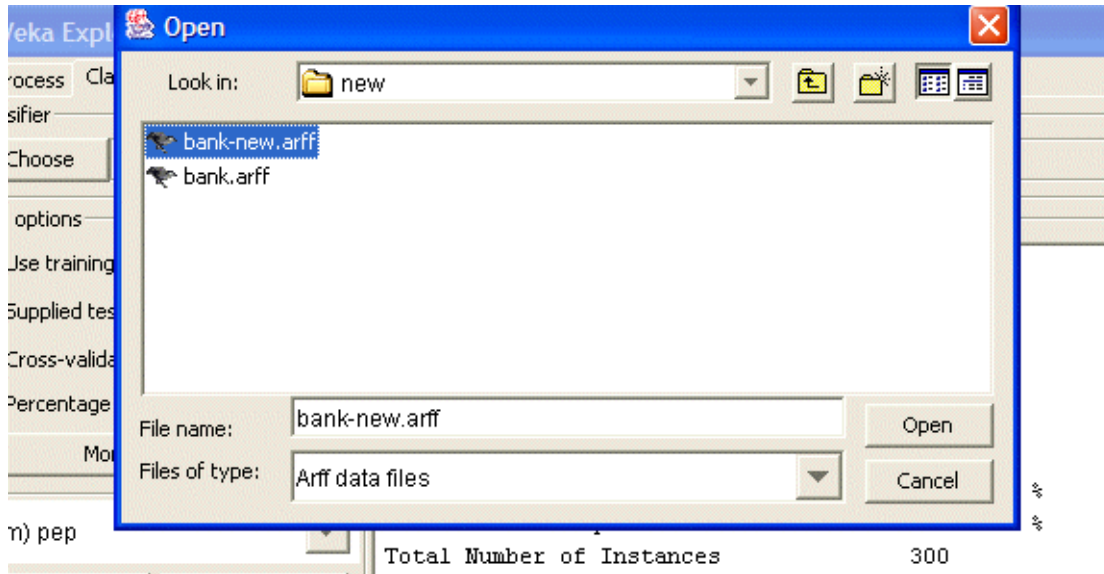
=== Detailed Accuracy By Class ===

TP Rate	FP Rate	Precision	Recall	F-Measure	Class
0.536	0.185	0.712	0.536	0.612	YES
0.815	0.464	0.673	0.815	0.737	NO

=== Confusion Matrix ===

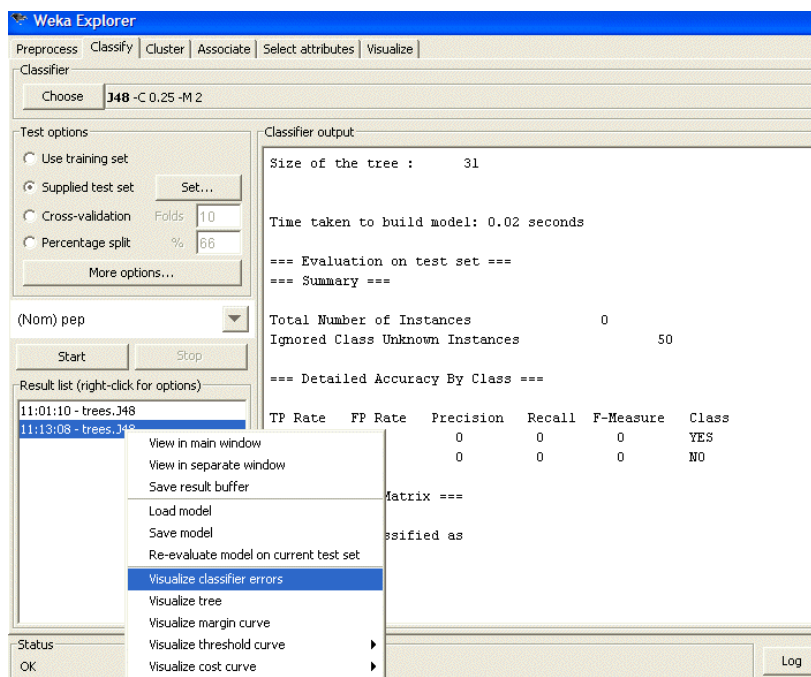
a	b	<-- classified as	
74	64	a = YES	
30	132	b = NO	

Figure 27-a



[Figure 27-b](#)

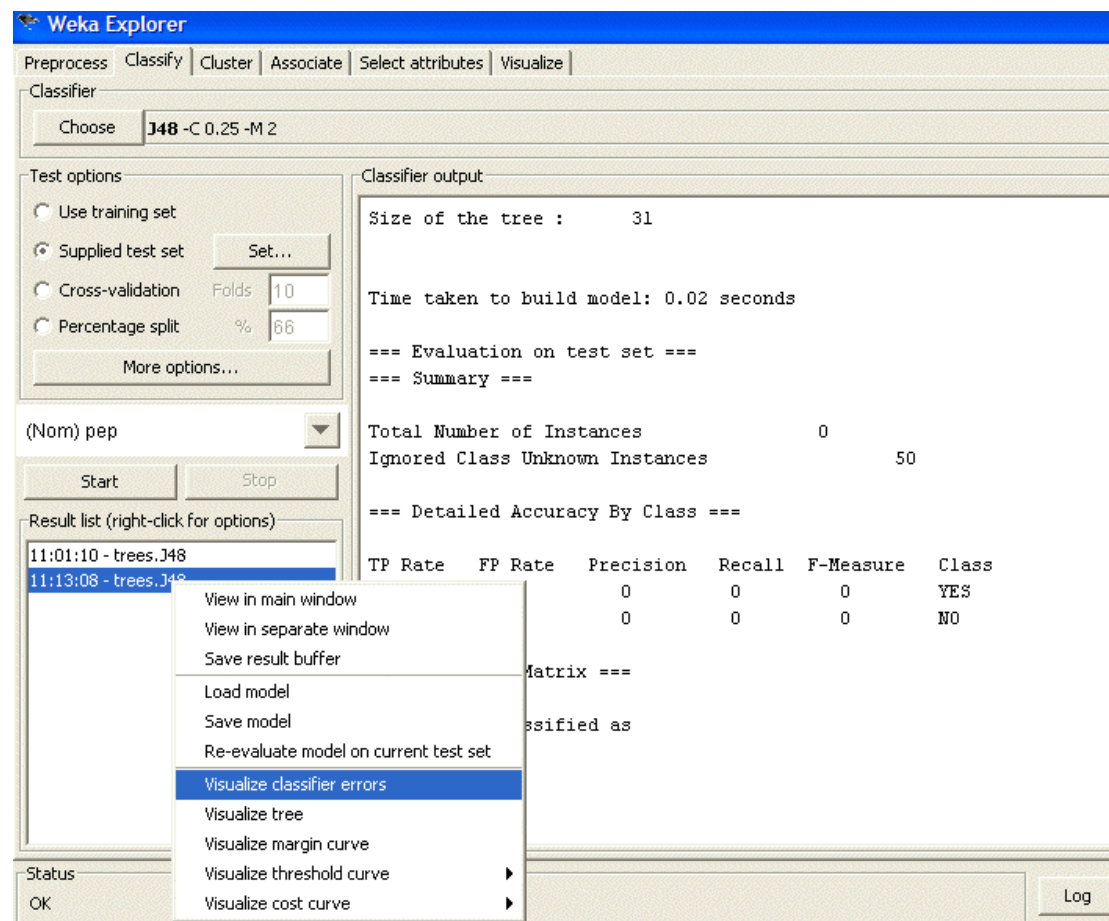
In this case, we open the file "bank-new.arff" and upon returning to the main window, we click the "start" button. This, once again generates the models from our training data, but this time it applies the model to the new unclassified instances in the "bank-new.arff" file in order to predict the value of "pep" attribute. The result is depicted in Figure 28. Note that the summary of the results in the right panel does not show any statistics. This is because in our test instances the value of the class attribute ("pep") was left as "?", thus WEKA has no actual values to which it can compare the predicted values of new instances.



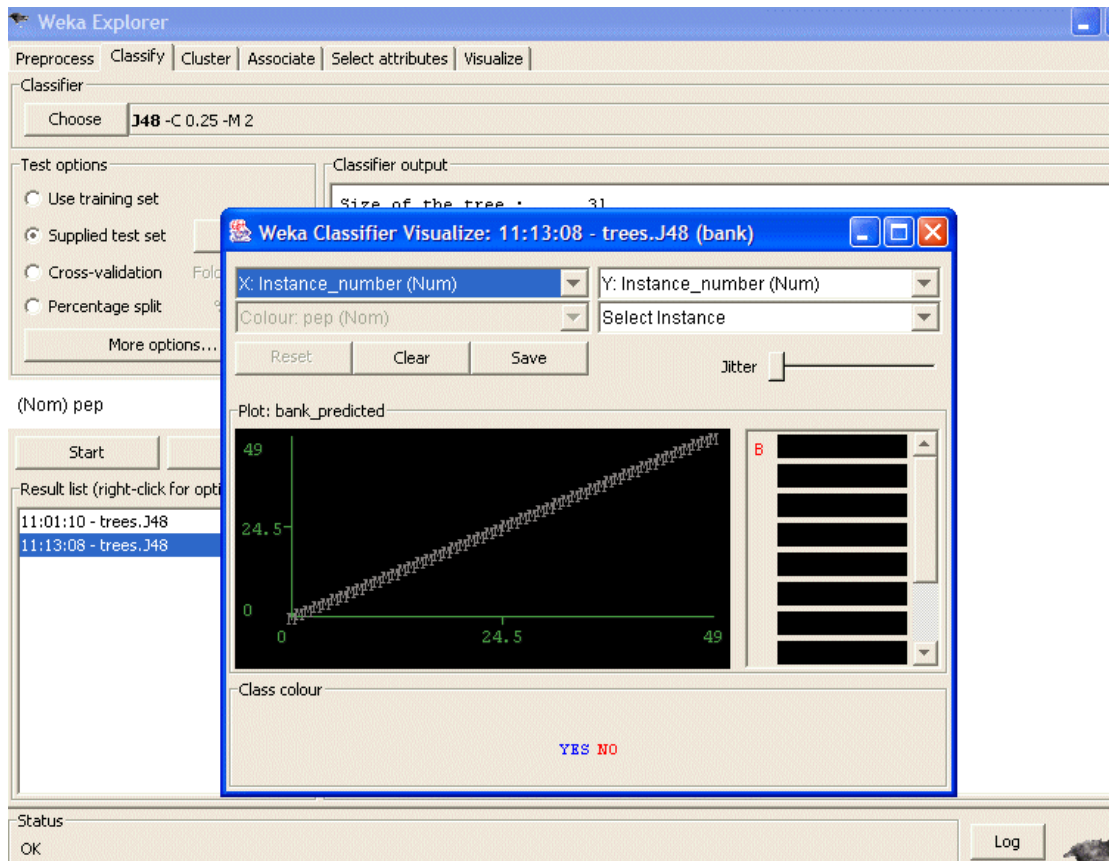
[Figure 28](#)

Of course, in this example we are interested in knowing how our model managed to classify the new instances. To do so we need to create a file containing all the new instances along with their predicted class value resulting from the application of the model. Doing this is much simpler using the command line version of WEKA classifier application. However, it is possible to do so in the GUI version using an "indirect" approach, as follows.

First, right-click the most recent result set in the left "Result list" panel. In the resulting pop-up window select the menu item "Visualize classifier errors". This brings up a separate window containing a two-dimensional graph. These steps and the resulting window are shown in Figures 28 and 29.

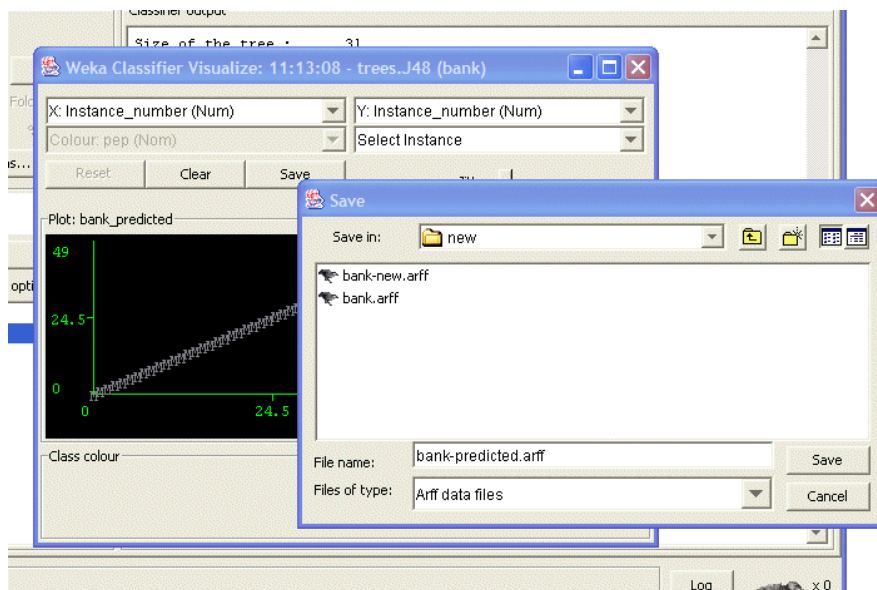


[Figure 28](#)



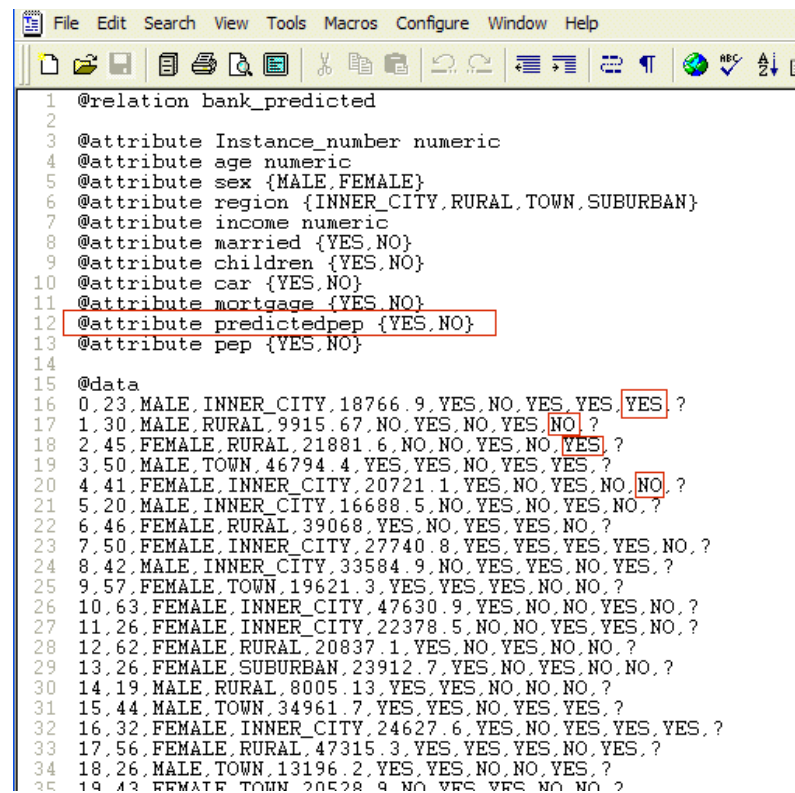
[Figure 29](#)

For now, we are not interested in what this graph represents. Rather, we would like to "save" the classification results from which the graph is generated. In the new window, we click on the "Save" button and save the result as the file: "bank-predicted.arff", as shown in Figure 30.



[Figure 30](#)

This file contains a copy of the new instances along with an additional column for the predicted value of "pep". The top portion of the file can be seen in Figure 31.



```

1 @relation bank_predicted
2
3 @attribute Instance_number numeric
4 @attribute age numeric
5 @attribute sex {MALE,FEMALE}
6 @attribute region {INNER_CITY,RURAL,TOWN,SUBURBAN}
7 @attribute income numeric
8 @attribute married {YES,NO}
9 @attribute children {YES,NO}
10 @attribute car {YES,NO}
11 @attribute mortgage {YES,NO}
12 @attribute predictedpep {YES,NO}
13 @attribute pep {YES,NO}
14
15 @data
16 0,23,MALE,INNER_CITY,18766,9,YES,NO,YES,YES,YES,?
17 1,30,MALE,RURAL,9915,67,NO,YES,NO,YES,NO,?
18 2,45,FEMALE,RURAL,21881,6,NO,NO,YES,NO,YES,?
19 3,50,MALE,TOWN,46794,4,YES,YES,NO,YES,YES,?
20 4,41,FEMALE,INNER_CITY,20721,1,YES,NO,YES,NO,NO,?
21 5,20,MALE,INNER_CITY,16688,5,NO,YES,NO,YES,NO,?
22 6,46,FEMALE,RURAL,39068,YES,NO,YES,YES,NO,?
23 7,50,FEMALE,INNER_CITY,27740,8,YES,YES,YES,YES,NO,?
24 8,42,MALE,INNER_CITY,33584,9,NO,YES,YES,NO,YES,?
25 9,57,FEMALE,TOWN,19621,3,YES,YES,YES,NO,NO,?
26 10,63,FEMALE,INNER_CITY,47630,9,YES,NO,NO,YES,NO,?
27 11,26,FEMALE,INNER_CITY,22378,5,NO,NO,YES,YES,NO,?
28 12,62,FEMALE,RURAL,20837,1,YES,NO,YES,NO,NO,?
29 13,26,FEMALE,SUBURBAN,23912,7,YES,NO,YES,NO,NO,?
30 14,19,MALE,RURAL,8005,13,YES,YES,NO,NO,NO,?
31 15,44,MALE,TOWN,34961,7,YES,YES,NO,YES,YES,?
32 16,32,FEMALE,INNER_CITY,24627,6,YES,NO,YES,YES,YES,?
33 17,56,FEMALE,RURAL,47315,3,YES,YES,YES,NO,YES,?
34 18,26,MALE,TOWN,13196,2,YES,YES,NO,NO,YES,?
35 19,43,FEMALE,TOWN,20528,9,NO,YES,YES,NO,NO,?

```

Figure 31

Note that two attributes have been added to the original new instances data: "Instance_number" and "predictedpep". These correspond to new columns in the data portion. The "predictedpep" value for each new instance is the last value before "?" which the actual "pep" class value. For example, the predicted value of the "pep" attribute for instance 0 is "YES" according to our model, while the predicted class value for instance 4 is "NO".