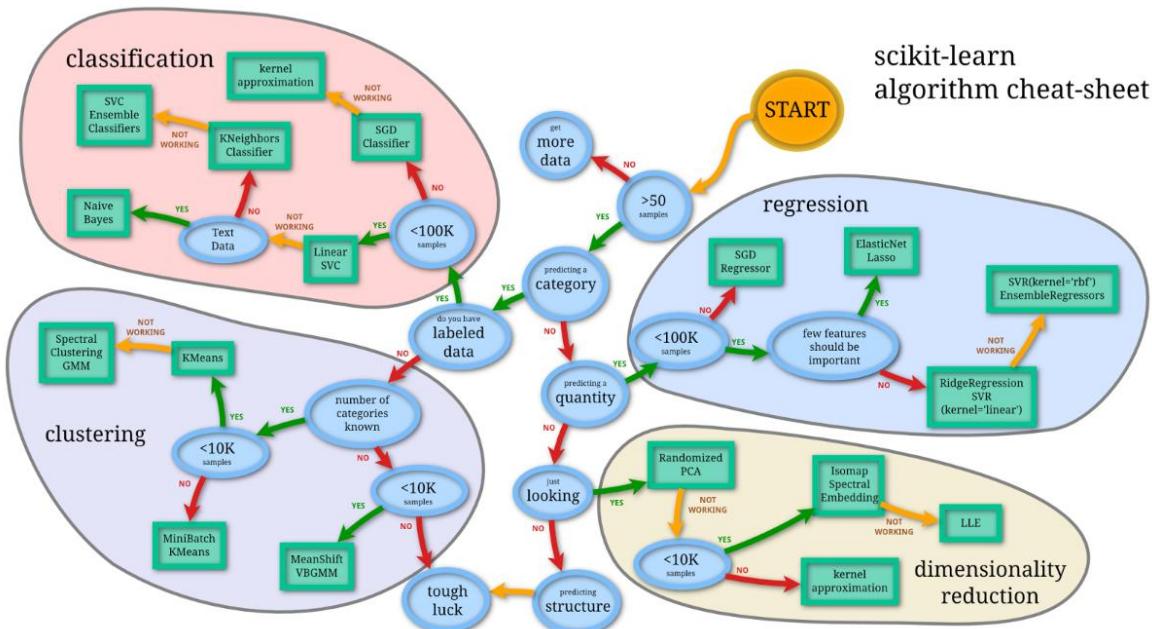


Overview: Data Mining Methods



WEKA Tutorial

- WEKA: A Machine Learning Toolkit
- The Explorer
 - Classification and Regression
 - Clustering
 - Association Rules
 - Attribute Selection
 - Data Visualization
- The Experimenter
- The Knowledge Flow GUI
- Conclusions

WEKA - Introduction

- Machine learning/data mining software written in Java (distributed under the GNU Public License)
- Used for research, education, and applications

- Main features:
 - Comprehensive set of data pre-processing tools, learning algorithms and evaluation methods
 - Graphical user interfaces (incl. data visualization)
 - Environment for comparing learning algorithms

5

Pre-processing the data

- Data can be imported from a file in various formats: ARFF, CSV, C4.5, binary
- Data can also be read from a URL or from an SQL database (using JDBC)
- Pre-processing tools in WEKA are called “filters”
- WEKA contains filters for:
 - Discretization, normalization, resampling, attribute selection, transforming and combining attributes, ...

6

WEKA with “flat” files

```
@relation heart-disease-simplified  
  
@attribute age numeric  
@attribute sex { female, male}  
@attribute chest_pain_type { typ_angina, asympt, non_anginal, atyp_angina}  
@attribute cholesterol numeric  
@attribute exercise_induced_angina { no, yes}  
@attribute class { present, not_present}  
  
@data  
63,male,typ_angina,233,no,not_present  
67,male,asympt,286,yes,present  
67,male,asympt,229,yes,present  
38,female,non_anginal,?,no,not_present  
...
```

Flat file in
ARFF format

7

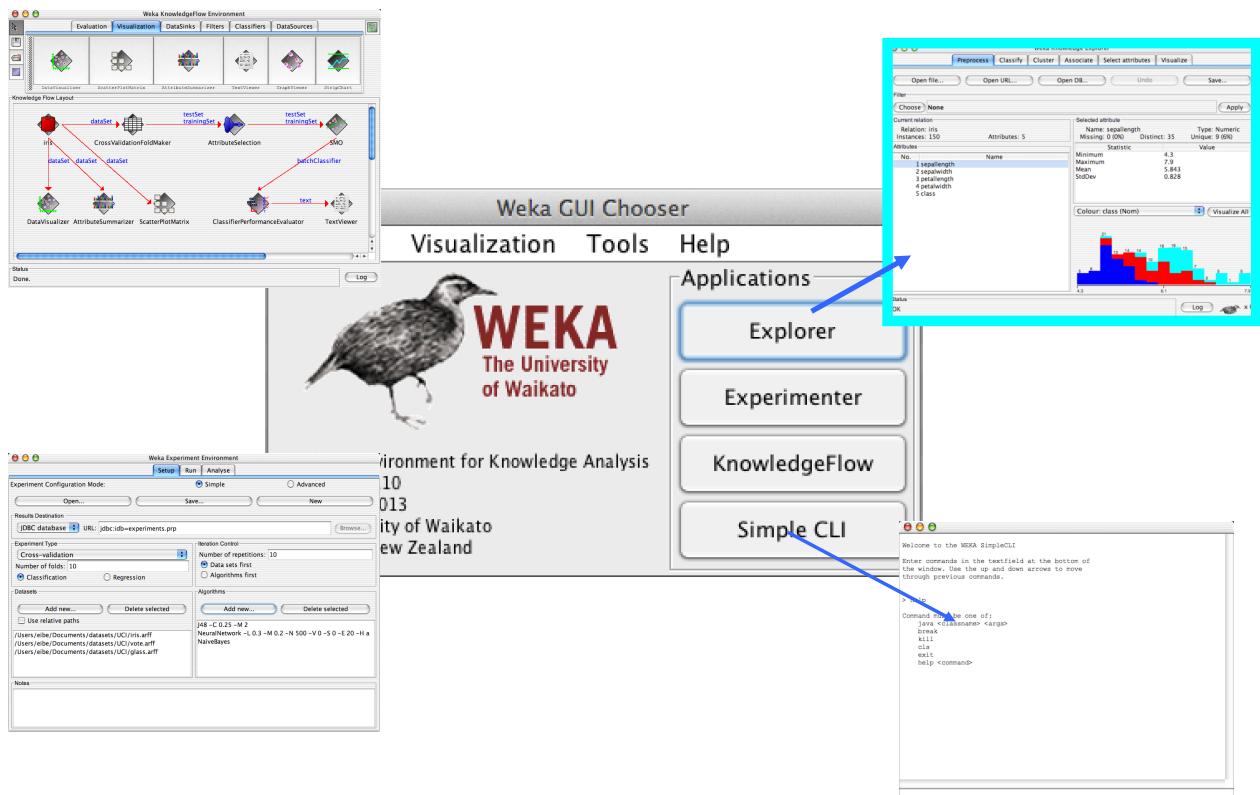
WEKA with “flat” files

```
@relation heart-disease-simplified  
  
@attribute age numeric  
@attribute sex { female, male}  
@attribute chest_pain_type { typ_angina, asympt, non_anginal, atyp_angina}  
@attribute cholesterol numeric  
@attribute exercise_induced_angina { no, yes}  
@attribute class { present, not_present}
```

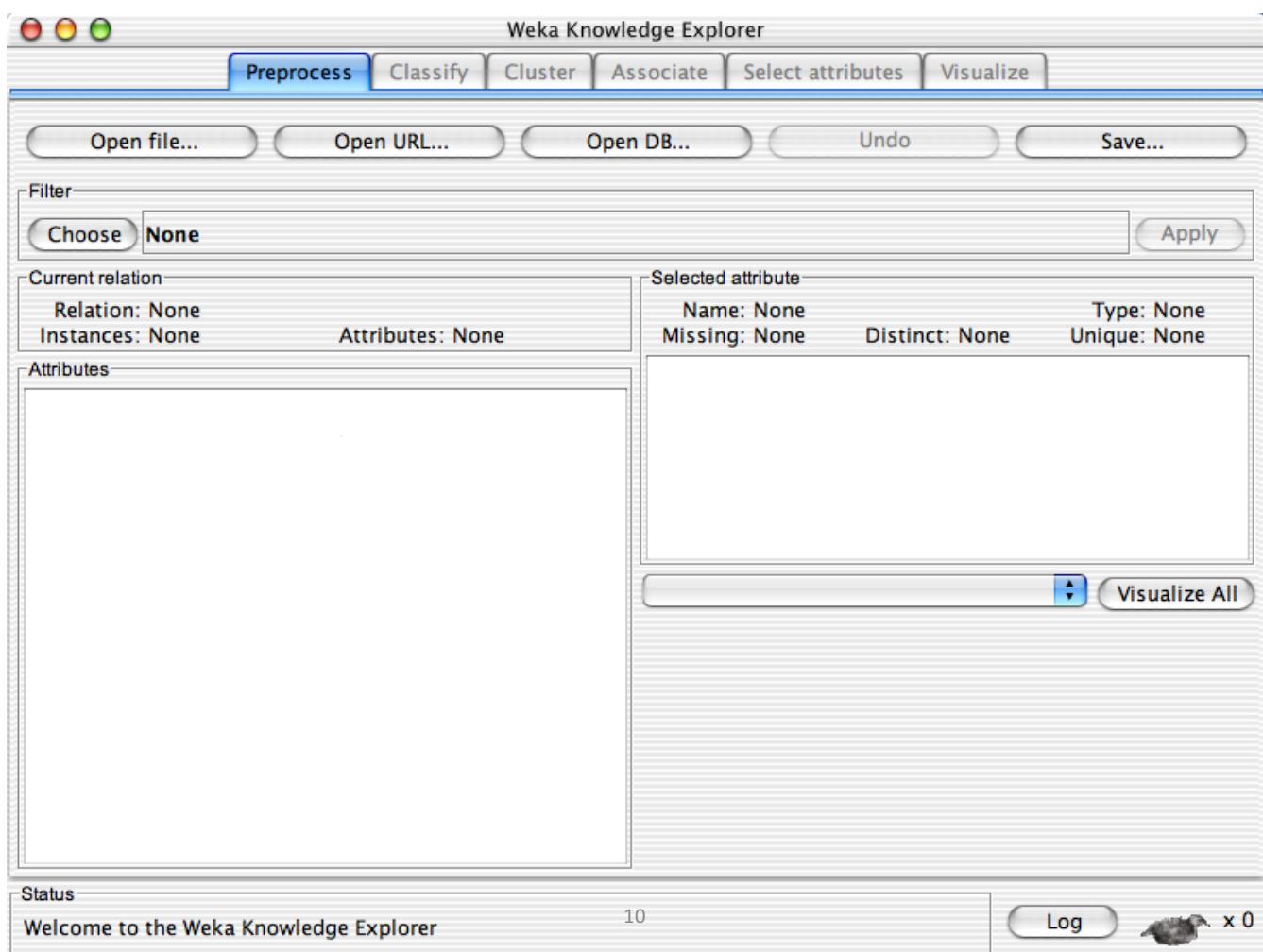
numeric attribute
nominal attribute

```
@data  
63,male,typ_angina,233,no,not_present  
67,male,asympt,286,yes,present  
67,male,asympt,229,yes,present  
38,female,non_anginal,?,no,not_present  
...
```

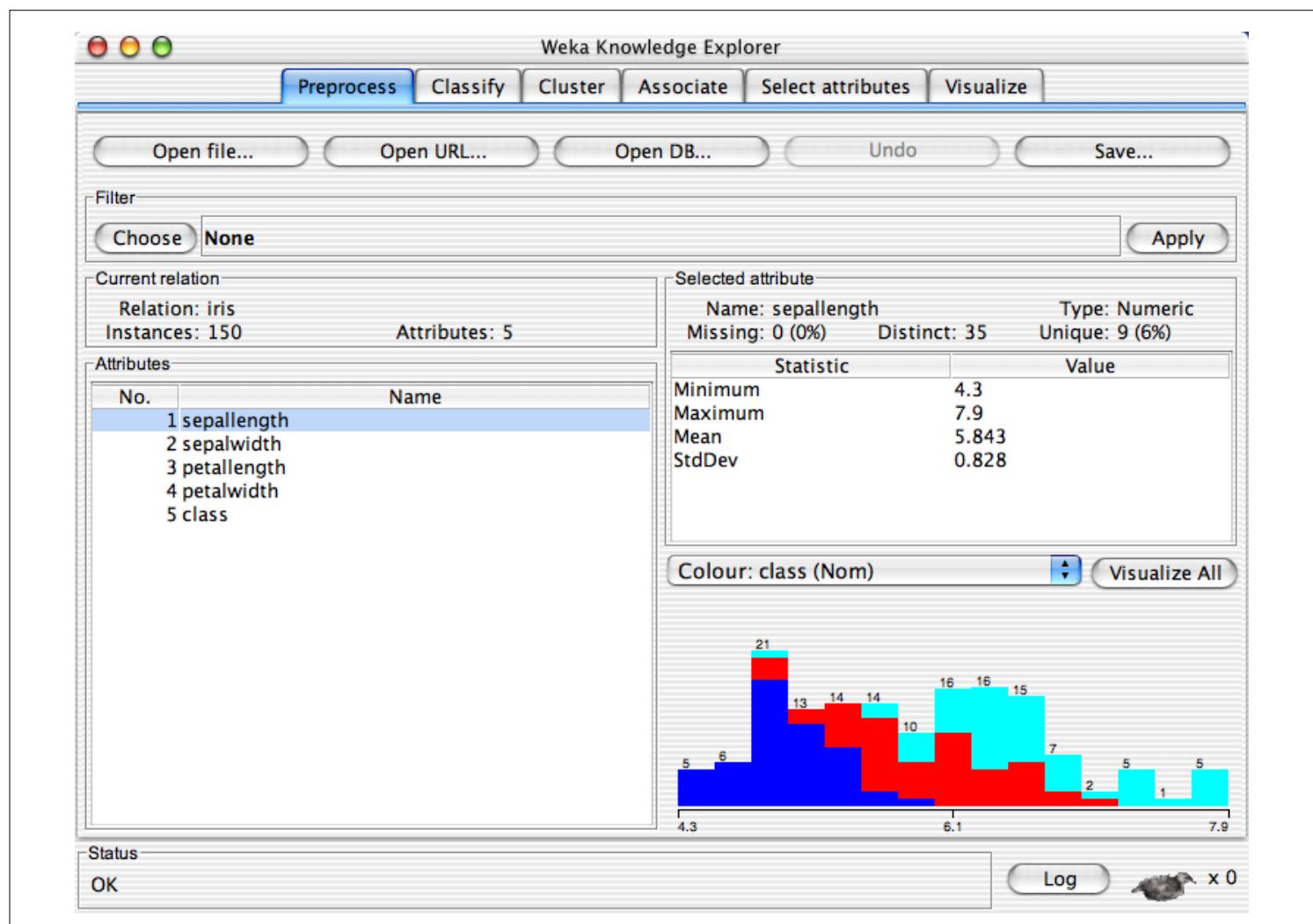
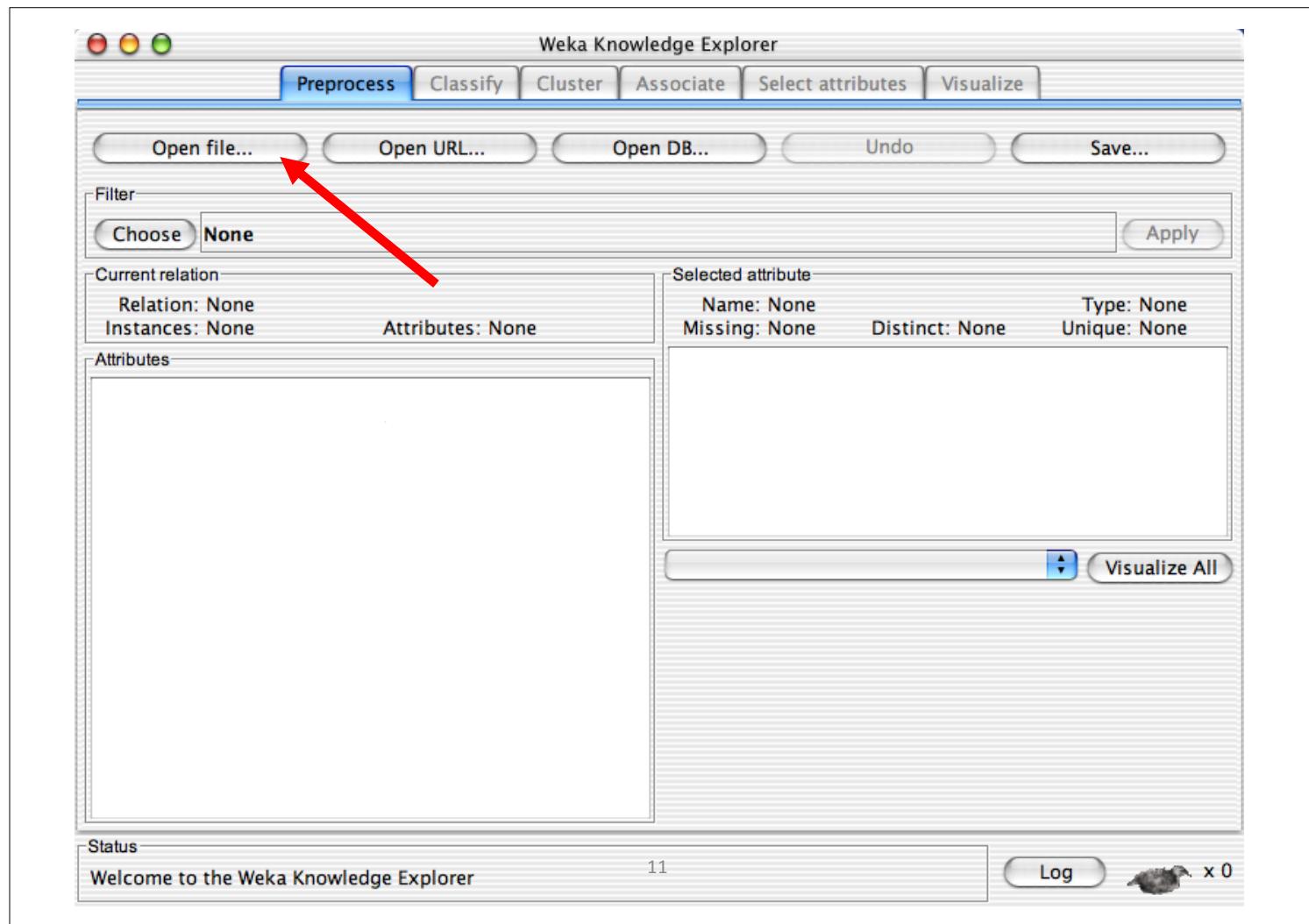
8



9



10



Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose None Apply

Current relation
Relation: iris Instances: 150 Attributes: 5

Attributes

No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

Selected attribute
Name: sepallength Type: Numeric
Missing: 0 (0%) Distinct: 35 Unique: 9 (6%)

Statistic	Value
Minimum	4.3
Maximum	7.9
Mean	5.843
StdDev	0.828

Colour: class (Nom) Visualize All

Status OK Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose None Apply

Current relation
Relation: iris Instances: 150 Attributes: 5

Attributes

No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

Selected attribute
Name: class Type: Nominal
Missing: 0 (0%) Distinct: 3 Unique: 0 (0%)

Label	Count
Iris-setosa	50
Iris-versicolor	50
Iris-virginica	50

Colour: class (Nom) Visualize All

Status OK Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose None Apply

Current relation
Relation: iris Instances: 150 Attributes: 5

Attributes

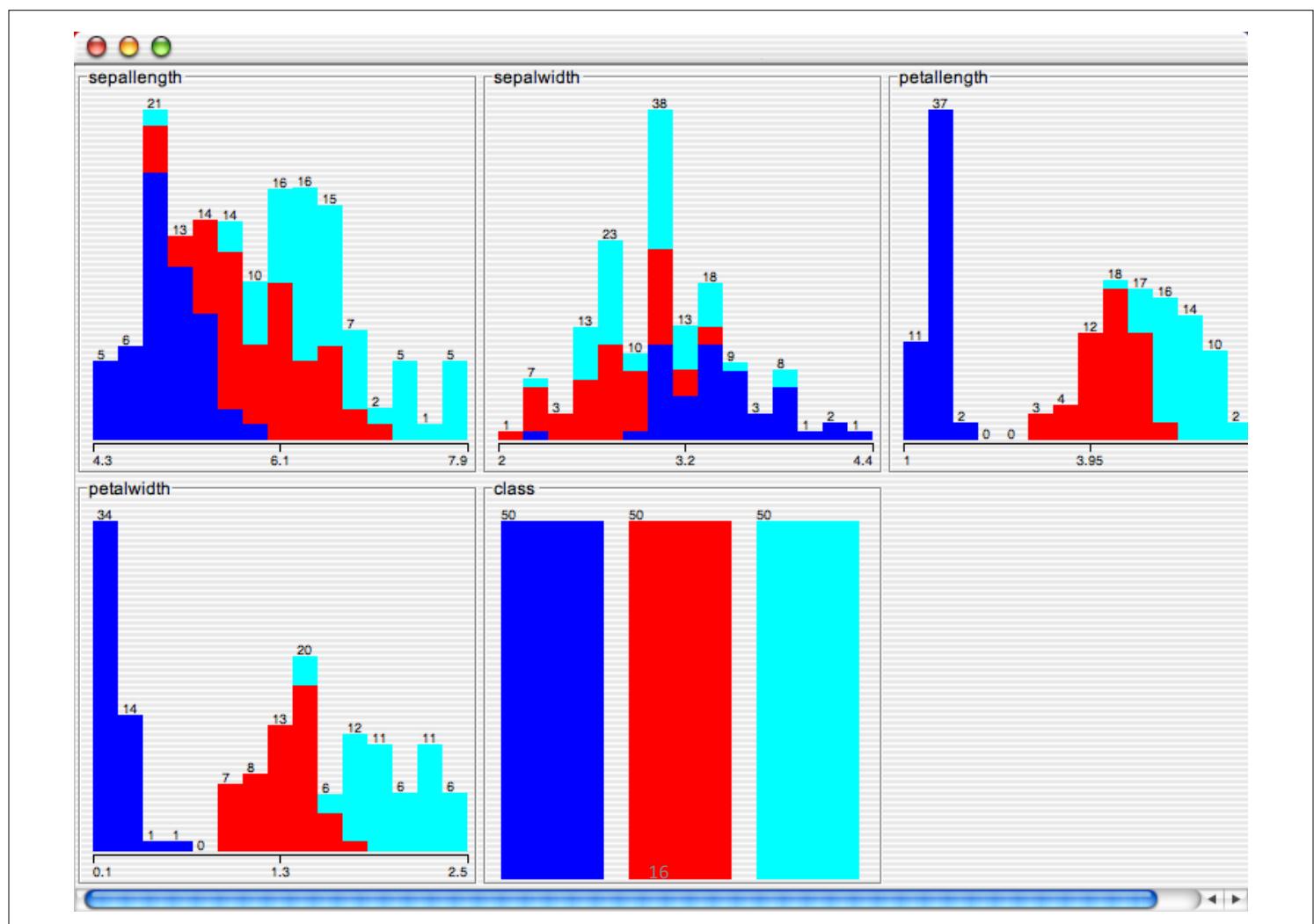
No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

Selected attribute
Name: class Type: Nominal
Missing: 0 (0%) Distinct: 3 Unique: 0 (0%)

Label	Count
Iris-setosa	50
Iris-versicolor	50
Iris-virginica	50

Colour: class (Nom) Visualize All

Status OK Log x 0



Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose None Apply

Current relation
Relation: iris Instances: 150 Attributes: 5

Attributes

No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

Selected attribute
Name: petallength Type: Numeric
Missing: 0 (0%) Distinct: 43 Unique: 10 (7%)

Statistic	Value
Minimum	1
Maximum	6.9
Mean	3.759
StdDev	1.764

Colour: class (Nom) Visualize All

Status OK Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose **None** Apply

Current relation
Relation: iris Instances: 150 Attributes: 5

Attributes

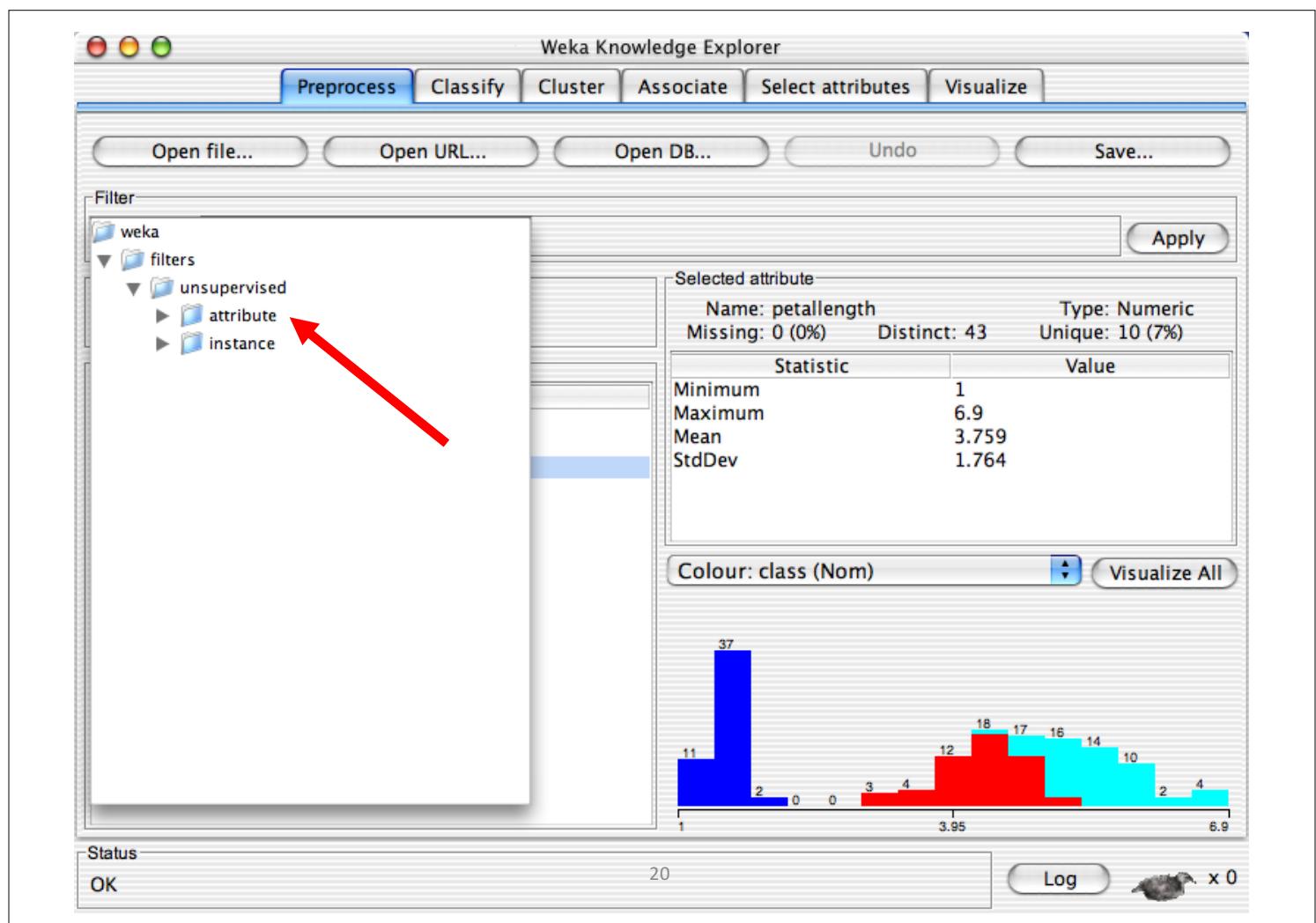
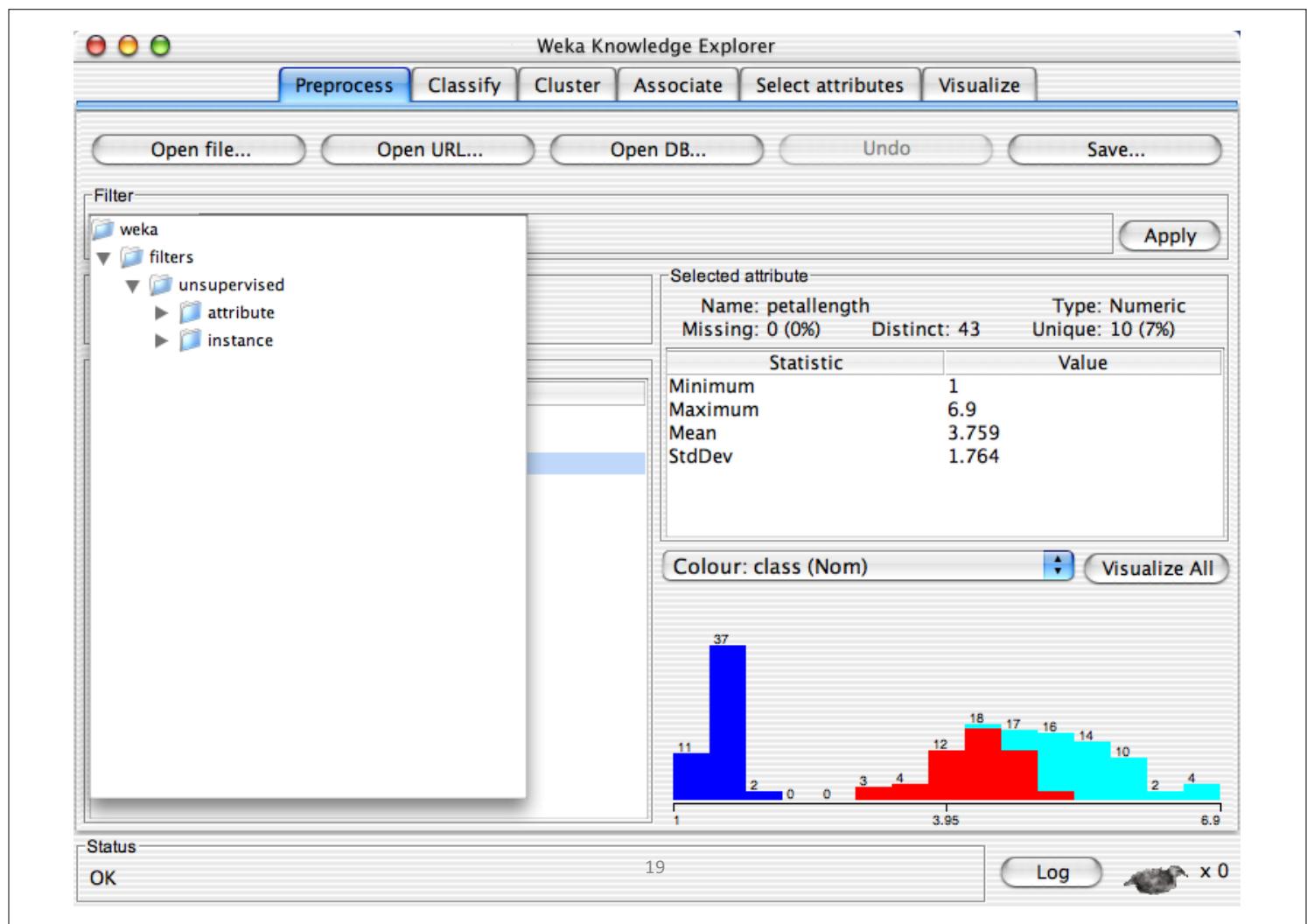
No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

Selected attribute
Name: petallength Type: Numeric
Missing: 0 (0%) Distinct: 43 Unique: 10 (7%)

Statistic	Value
Minimum	1
Maximum	6.9
Mean	3.759
StdDev	1.764

Colour: class (Nom) Visualize All

Status OK Log x 0



Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter

- weka
- filters
- unsupervised
 - attribute
 - Add
 - AddCluster
 - AddExpression
 - AddNoise
 - Copy
 - Discretize**
 - FirstOrder
 - MakeIndicator
 - MergeTwoValues
 - NominalToBinary
 - Normalize
 - NumericToBinary
 - NumericTransform
 - Obfuscate
 - PKIDiscretize
 - Remove
 - RemoveType

Selected attribute

Name: petallength Type: Numeric
Missing: 0 (0%) Distinct: 43 Unique: 10 (7%)

Statistic	Value
Minimum	1
Maximum	6.9
Mean	3.759
StdDev	1.764

Colour: class (Nom) Visualize All

Status: OK Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter

Choose **Discretize -B 10 -R first-last** Apply

Current relation

Relation: iris Instances: 150 Attributes: 5

Attributes

No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

Selected attribute

Name: petallength Type: Numeric
Missing: 0 (0%) Distinct: 43 Unique: 10 (7%)

Statistic	Value
Minimum	1
Maximum	6.9
Mean	3.759
StdDev	1.764

Colour: class (Nom) Visualize All

Status: OK Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose **Discretize -B 10 -R first-last** Apply

Current relation
Relation: iris Instances: 150 Attributes: 5

Attributes

No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

Selected attribute
Name: petallength Type: Numeric
Missing: 0 (0%) Distinct: 43 Unique: 10 (7%)

Statistic	Value
Minimum	1
Maximum	6.9
Mean	3.759
StdDev	1.764

Colour: class (Nom) Visualize All

Status OK Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose **Discretize -B 10 -R first-last** weka.gui.GenericObjectEditor Apply

Current relation
Relation: iris Instances: 150 Attributes: 5

Attributes

No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

weka.filters.unsupervised.attribute.Discretize

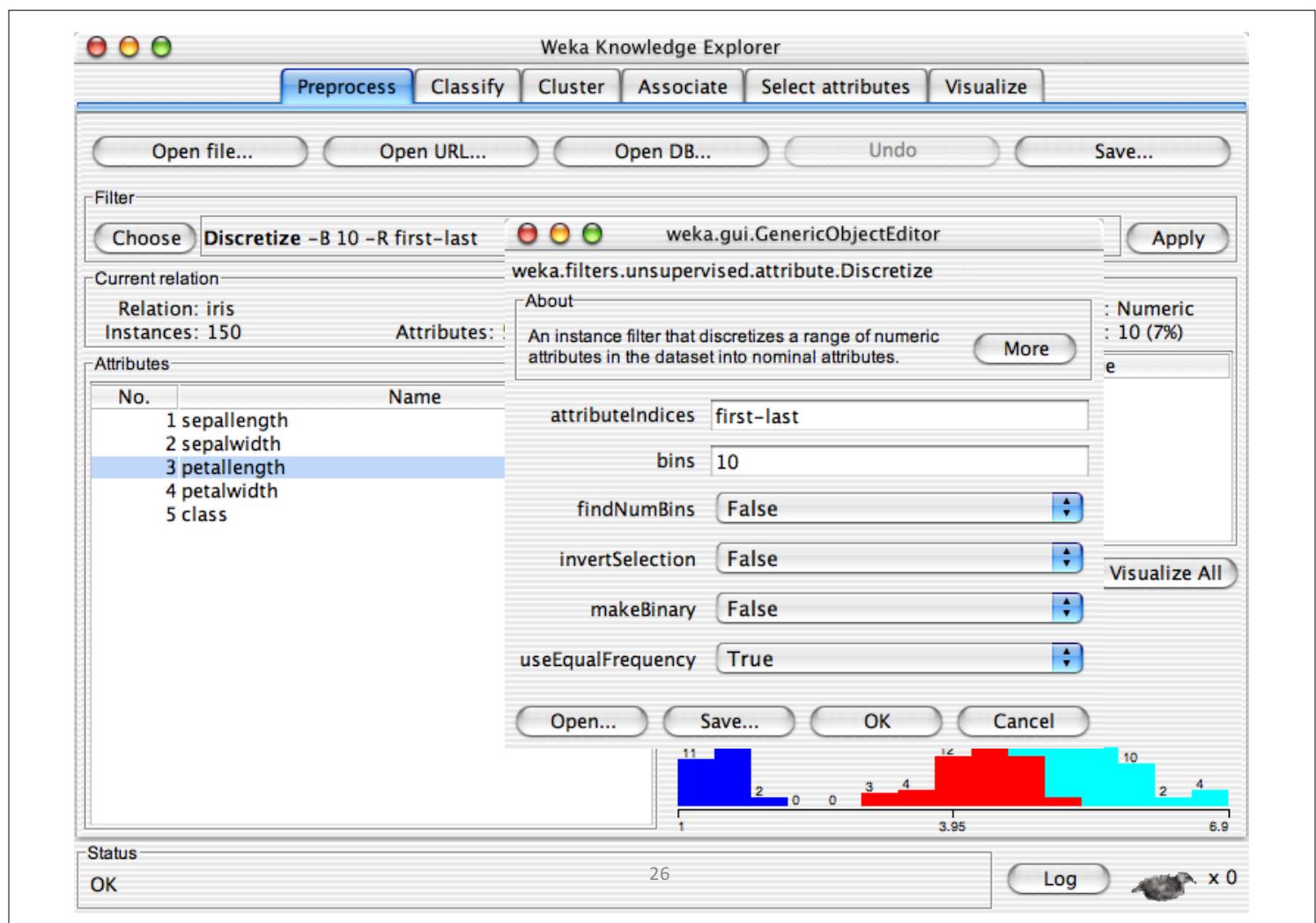
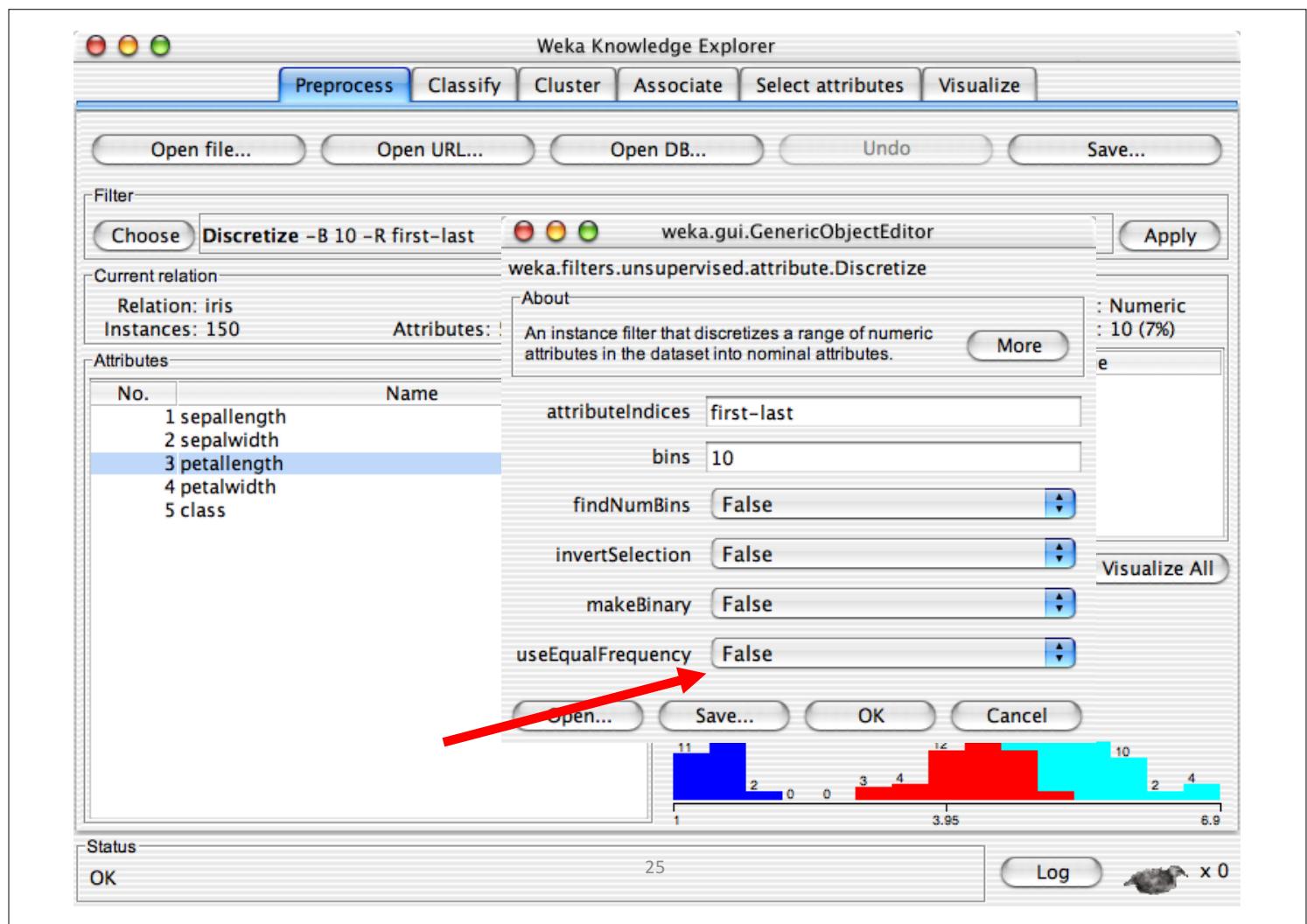
About An instance filter that discretizes a range of numeric attributes in the dataset into nominal attributes. More

attributeIndices first-last
bins 10
findNumBins False
invertSelection False
makeBinary False
useEqualFrequency False

Visualize All

Open... Save... OK Cancel

Status OK Log x 0



Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose Discretize -B 10 -R first-last weka.gui.GenericObjectEditor Apply

Current relation Relation: iris Instances: 150 Attributes: 5

Attributes

No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

About An instance filter that discretizes a range of numeric attributes in the dataset into nominal attributes. More

attributeIndices first-last
bins 10
findNumBins False
invertSelection False
makeBinary False
useEqualFrequency True

OK Cancel

Status OK Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose Discretize -F -B 10 -R first-last Apply

Current relation Relation: iris Instances: 150 Attributes: 5

Attributes

No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

Selected attribute Name: petallength Type: Numeric
Missing: 0 (0%) Distinct: 43 Unique: 10 (7%)

Statistic	Value
Minimum	1
Maximum	6.9
Mean	3.759
StdDev	1.764

Colour: class (Nom) Visualize All

Status OK Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose Discretize -F -B 10 -R first-last Apply

Current relation Relation: iris Instances: 150 Attributes: 5

Attributes

No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

Selected attribute Name: petallength Type: Numeric
Missing: 0 (0%) Distinct: 43 Unique: 10 (7%)

Statistic	Value
Minimum	1
Maximum	6.9
Mean	3.759
StdDev	1.764

Colour: class (Nom) Visualize All

Status OK Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Undo Save...

Filter Choose Discretize -F -B 10 -R first-last Apply

Current relation Relation: iris-weka.filters.unsupervised.attribute.Disc... Instances: 150 Attributes: 5

Attributes

No.	Name
1	sepallength
2	sepalwidth
3	petallength
4	petalwidth
5	class

Selected attribute Name: petallength Type: Nominal
Missing: 0 (0%) Distinct: 10 Unique: 0 (0%)

Label	Count
'(-inf-1.45]'	23
'(1.45-1.55]'	14
'(1.55-1.8]'	11
'(1.8-3.95]'	13
'(3.95-4.35]'	14
'(4.35-4.65]'	15
'(4.65-5.05]'	18

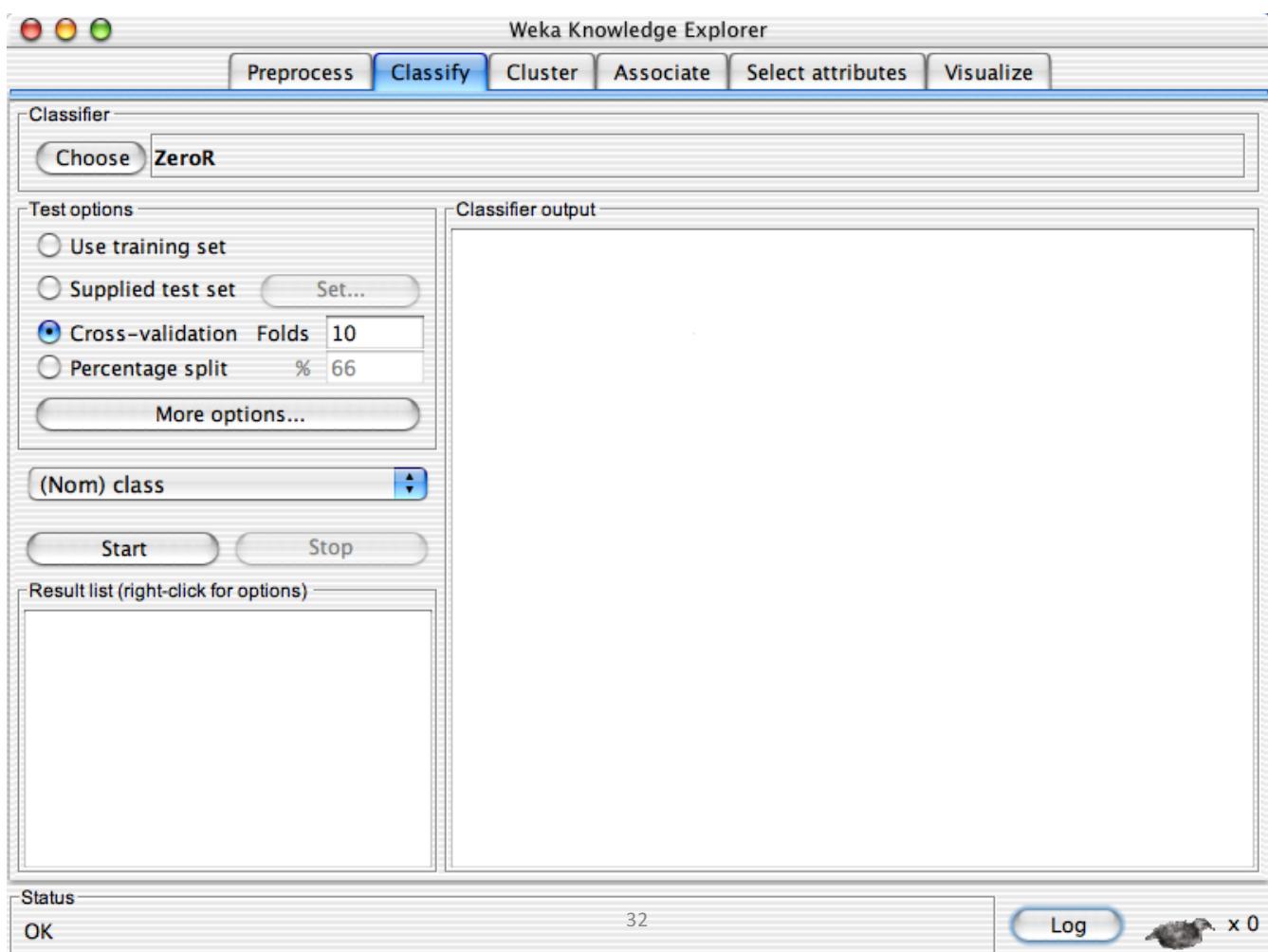
Colour: class (Nom) Visualize All

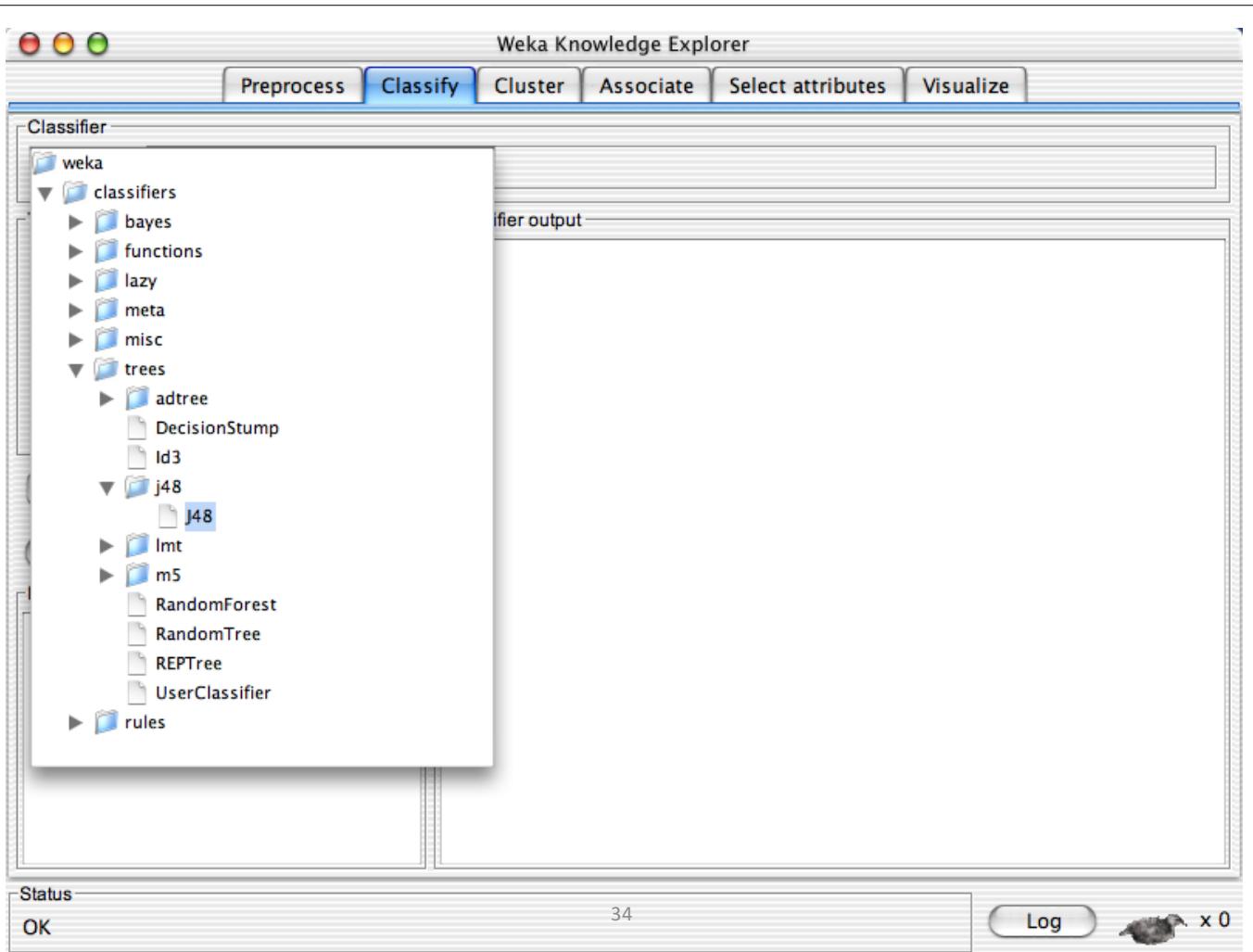
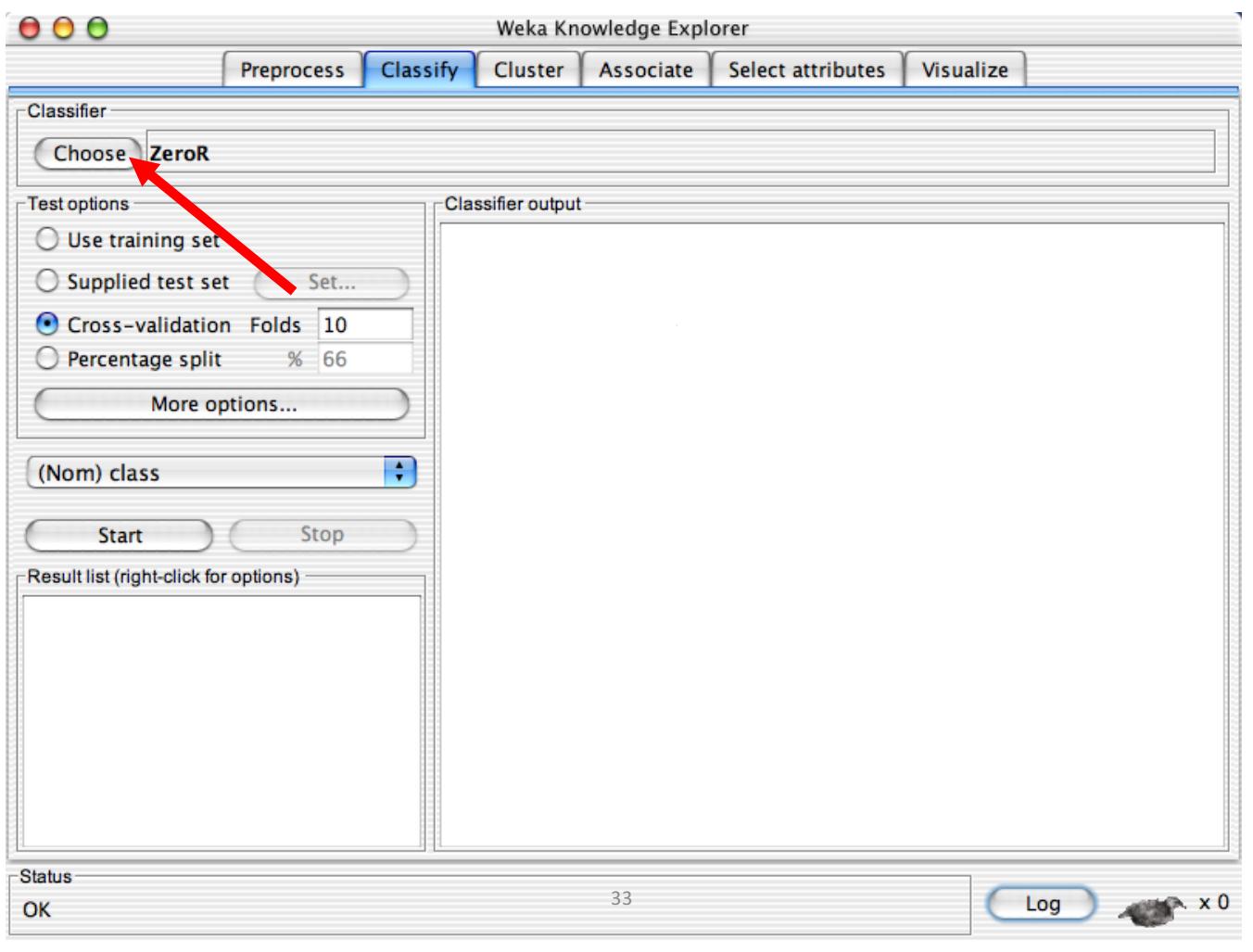
Status OK Log x 0

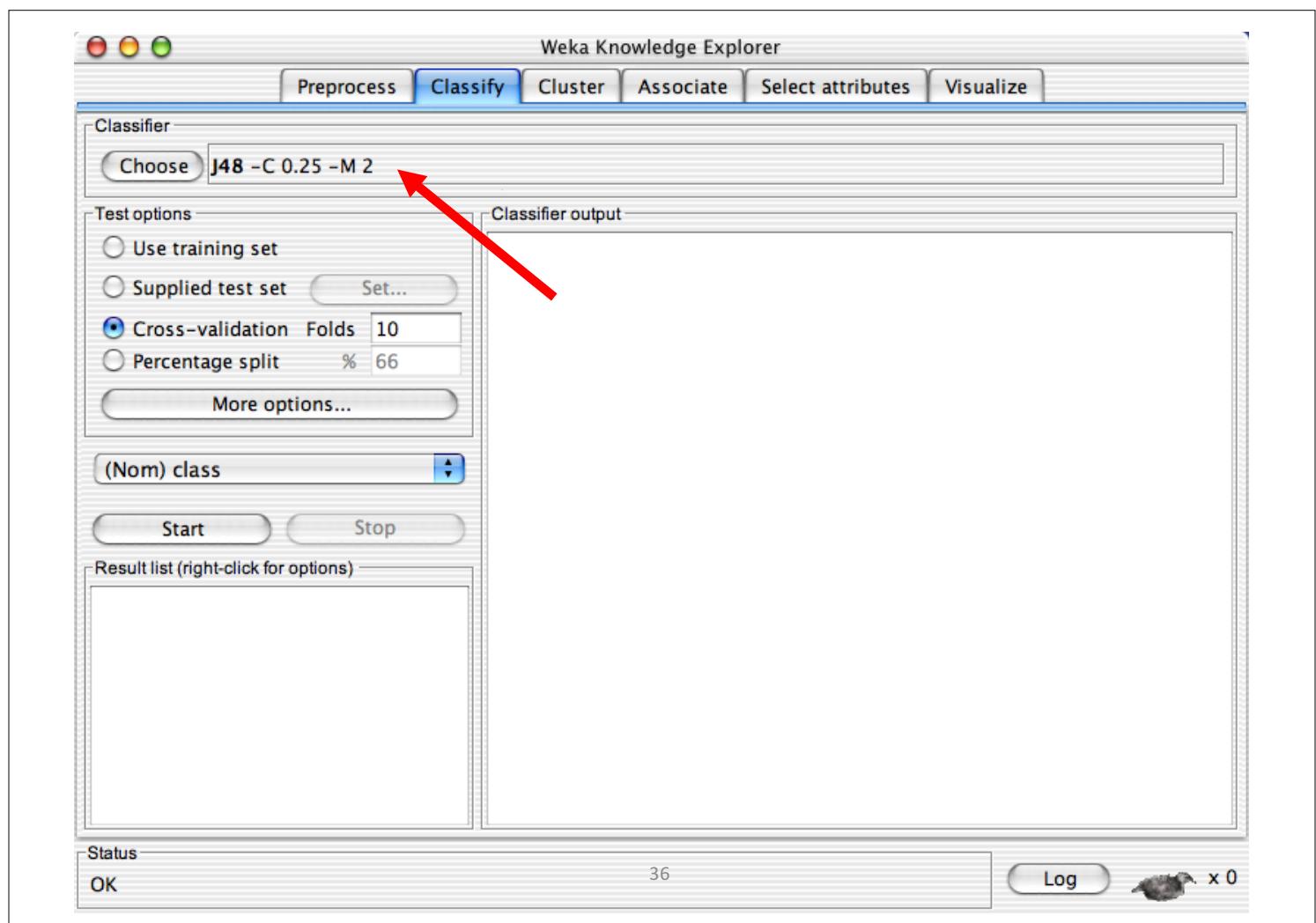
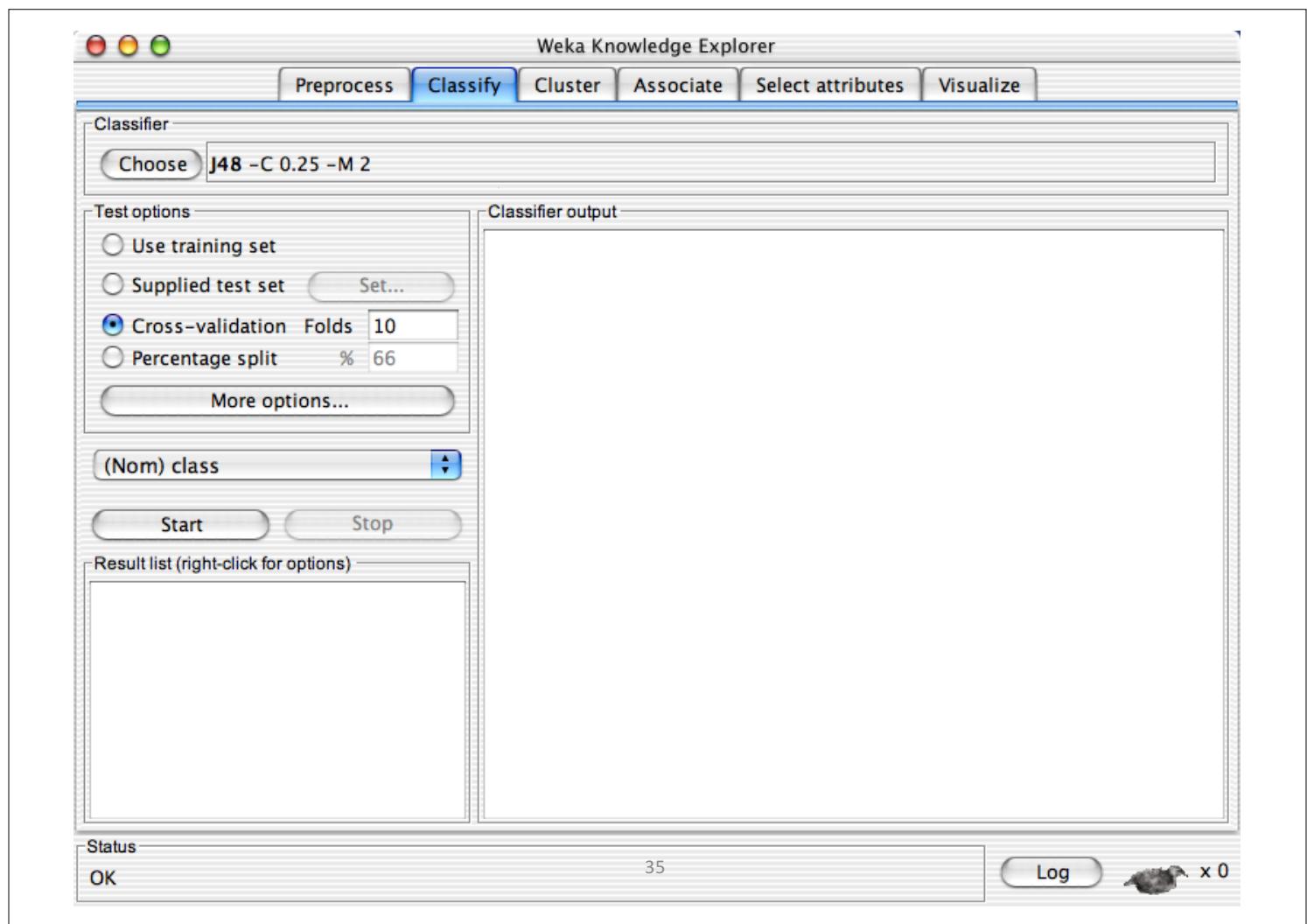
Building “Classifiers”

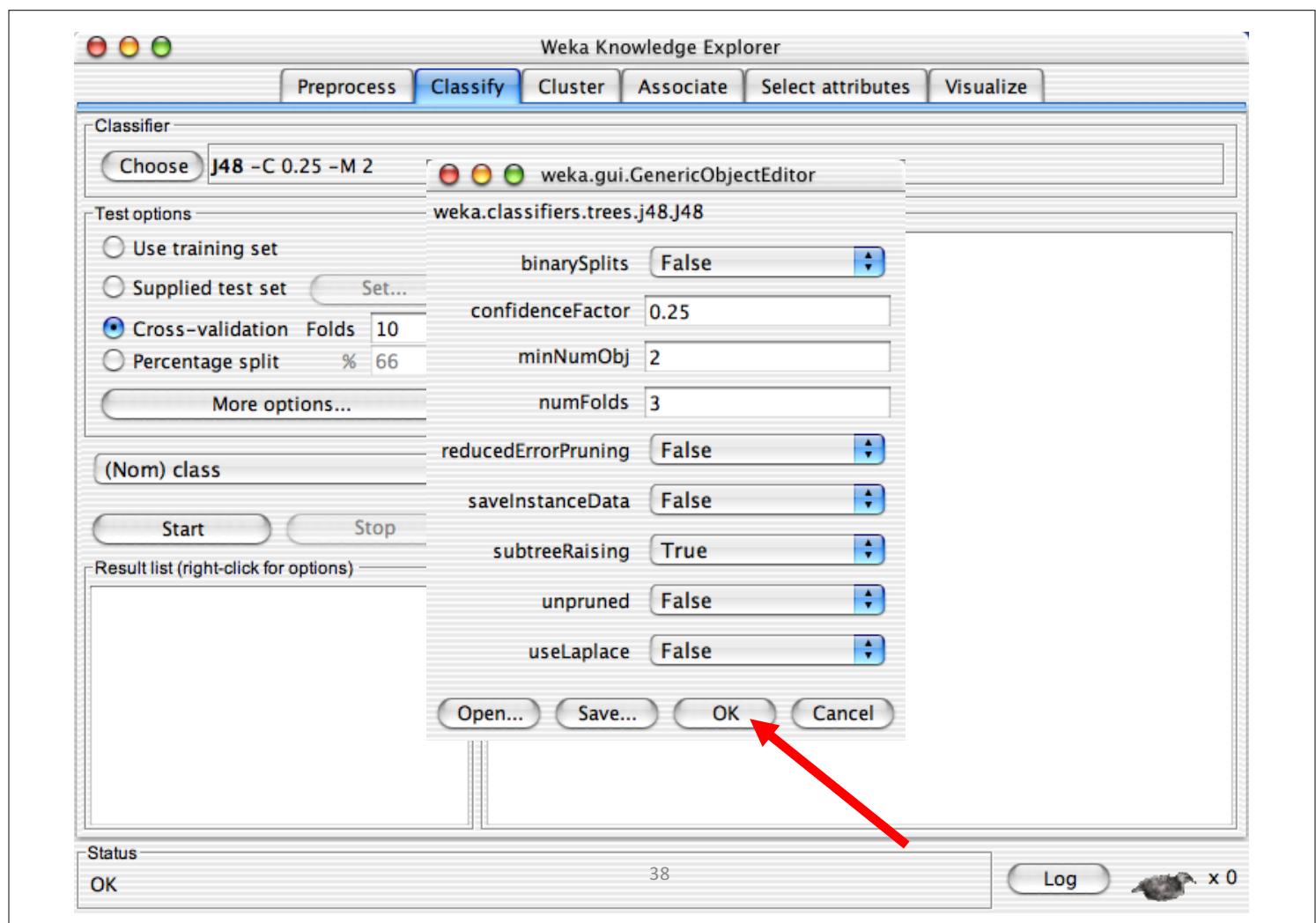
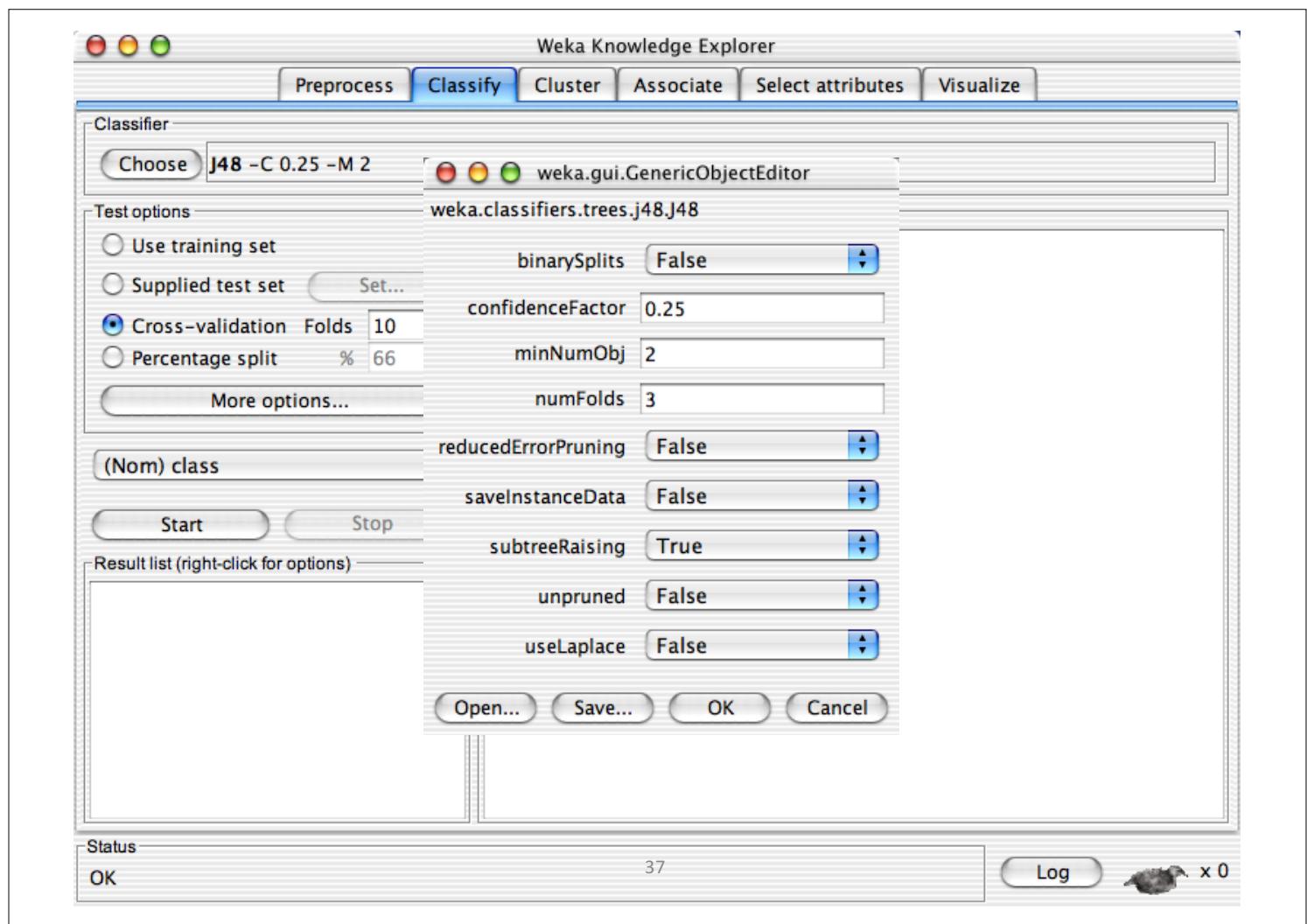
- Classifiers in WEKA are models for predicting nominal or numeric quantities
- Implemented learning schemes include:
 - Decision trees and lists, instance-based classifiers, support vector machines, multi-layer perceptrons, logistic regression, Bayes' nets, ...
- “Meta”-classifiers include:
 - Bagging, boosting, stacking, error-correcting output codes, locally weighted learning, ...

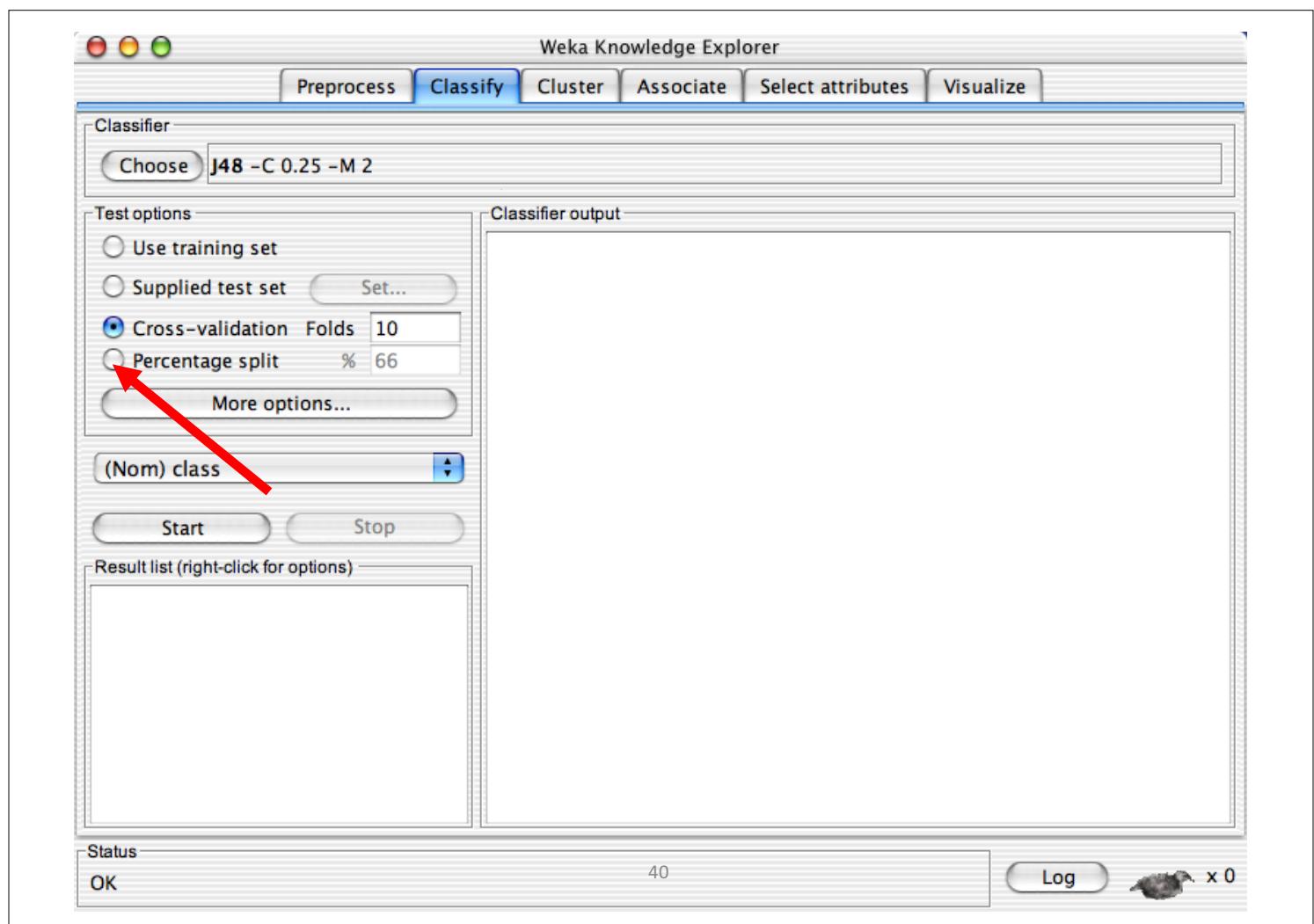
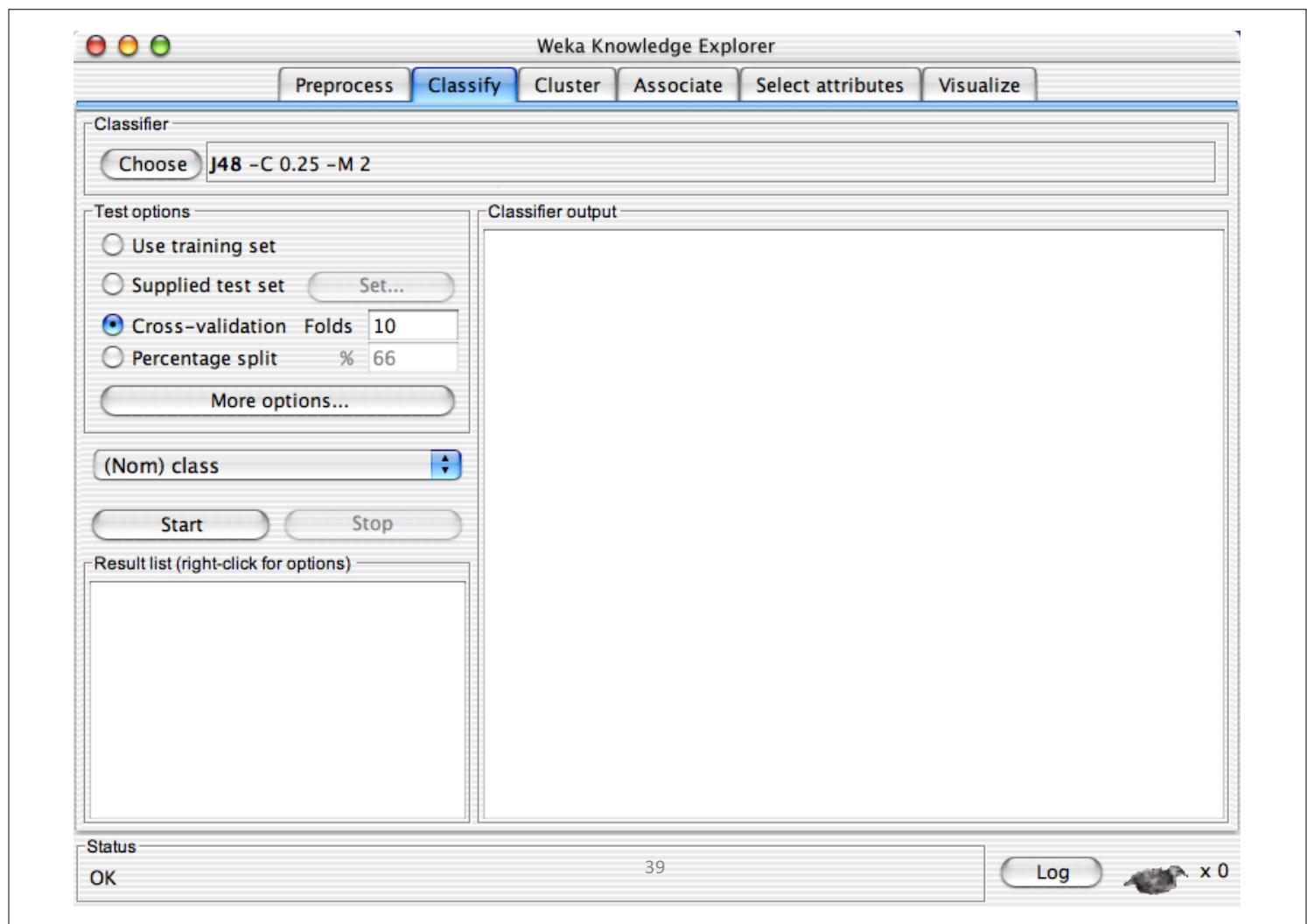
31

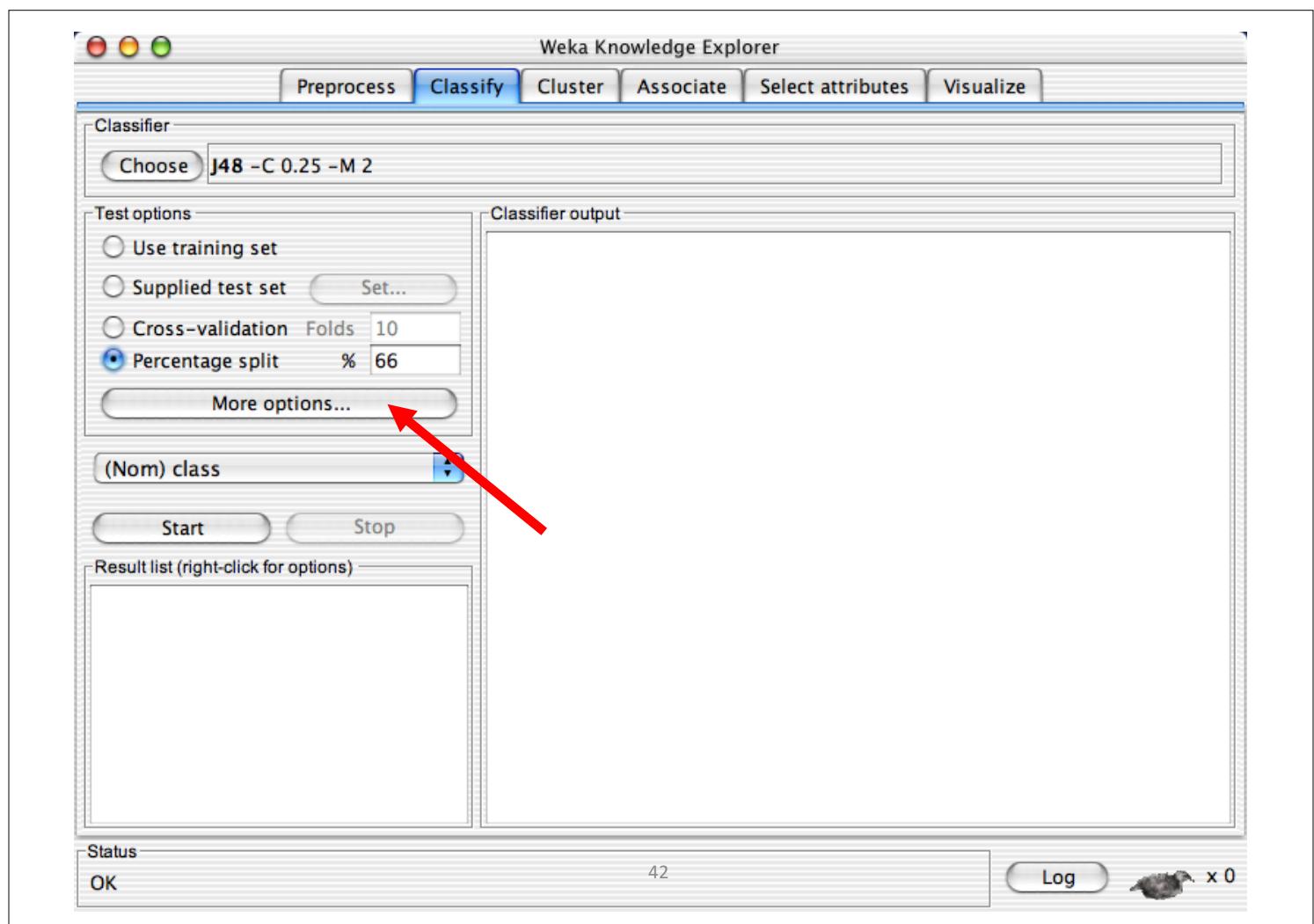
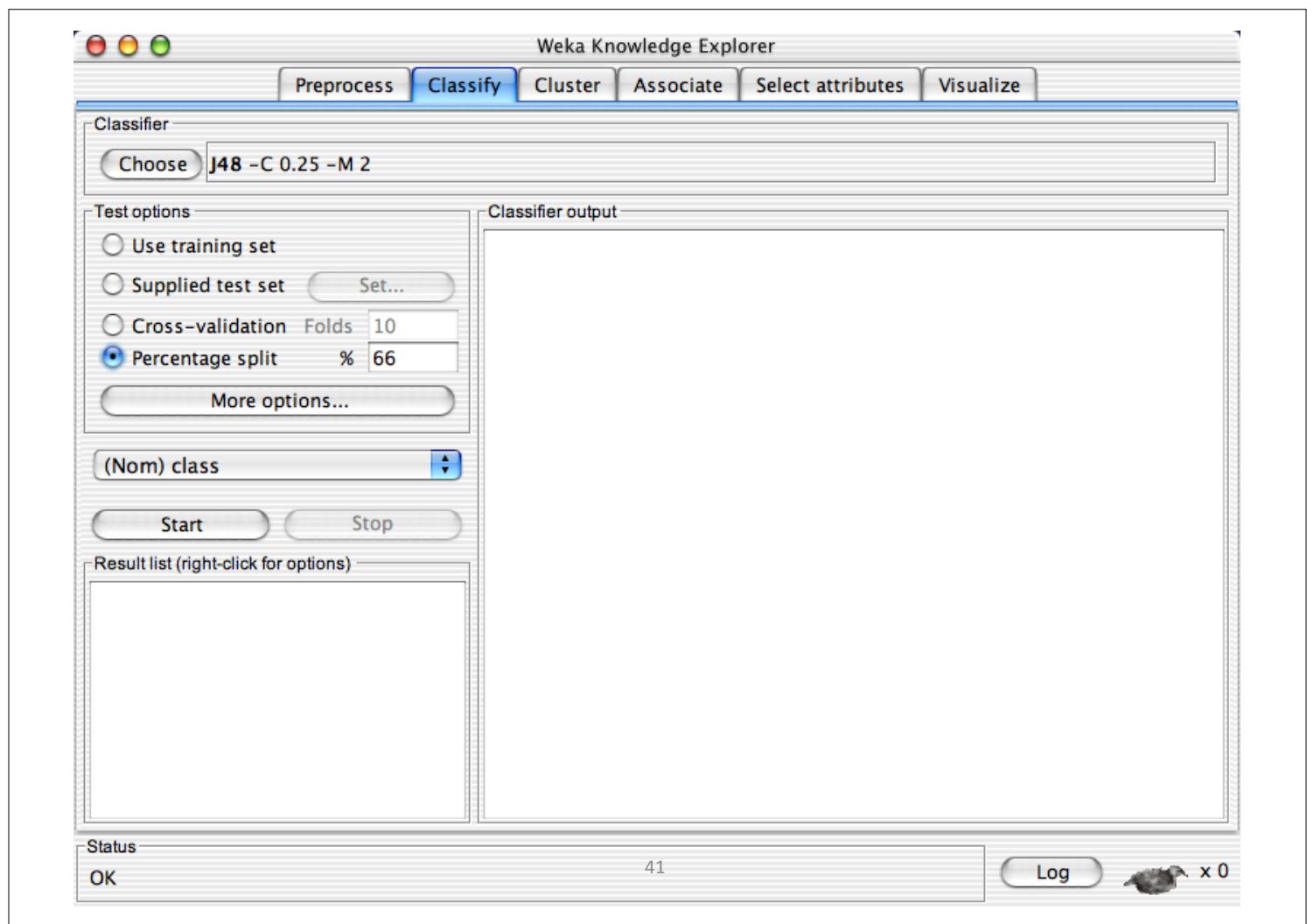


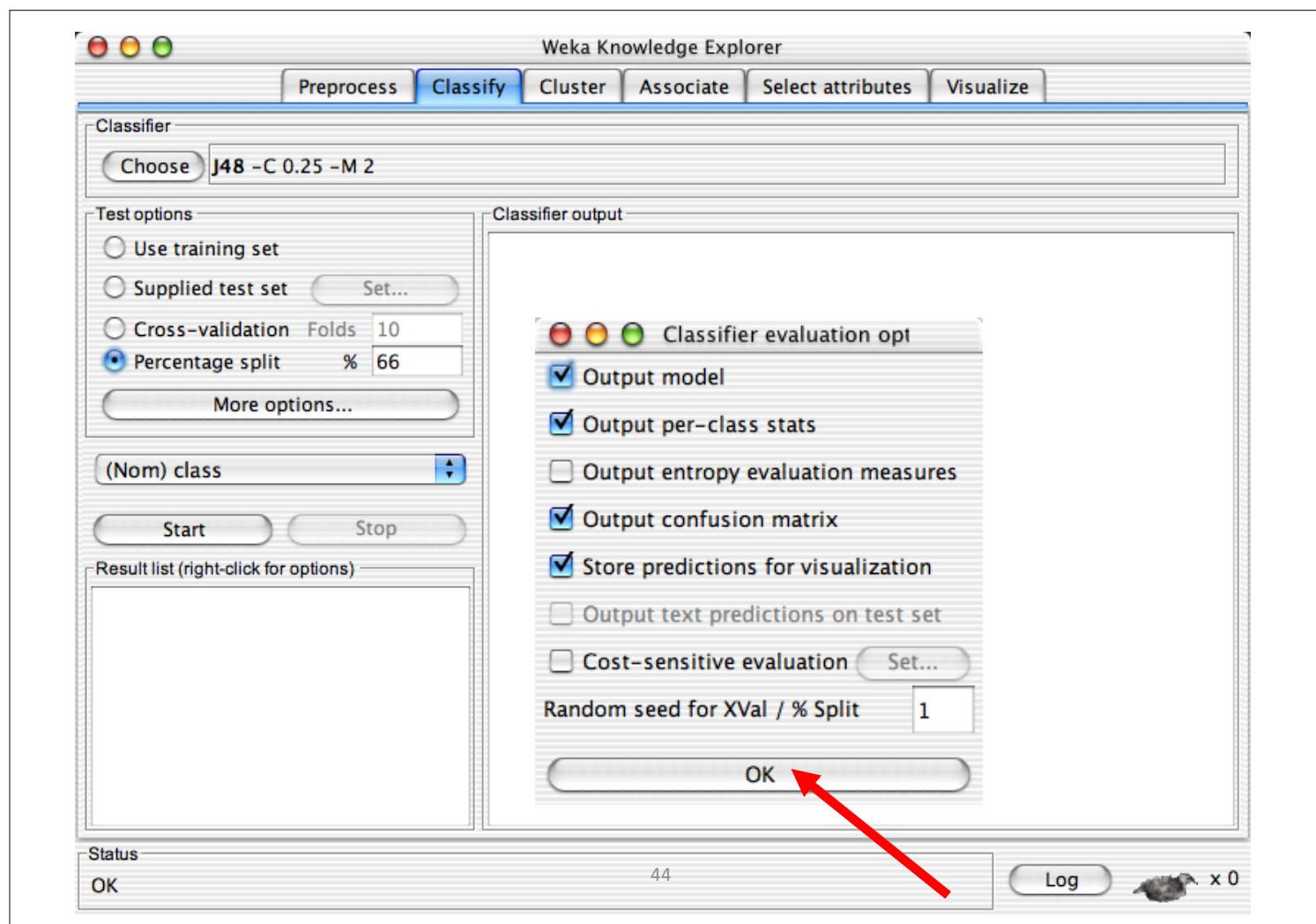
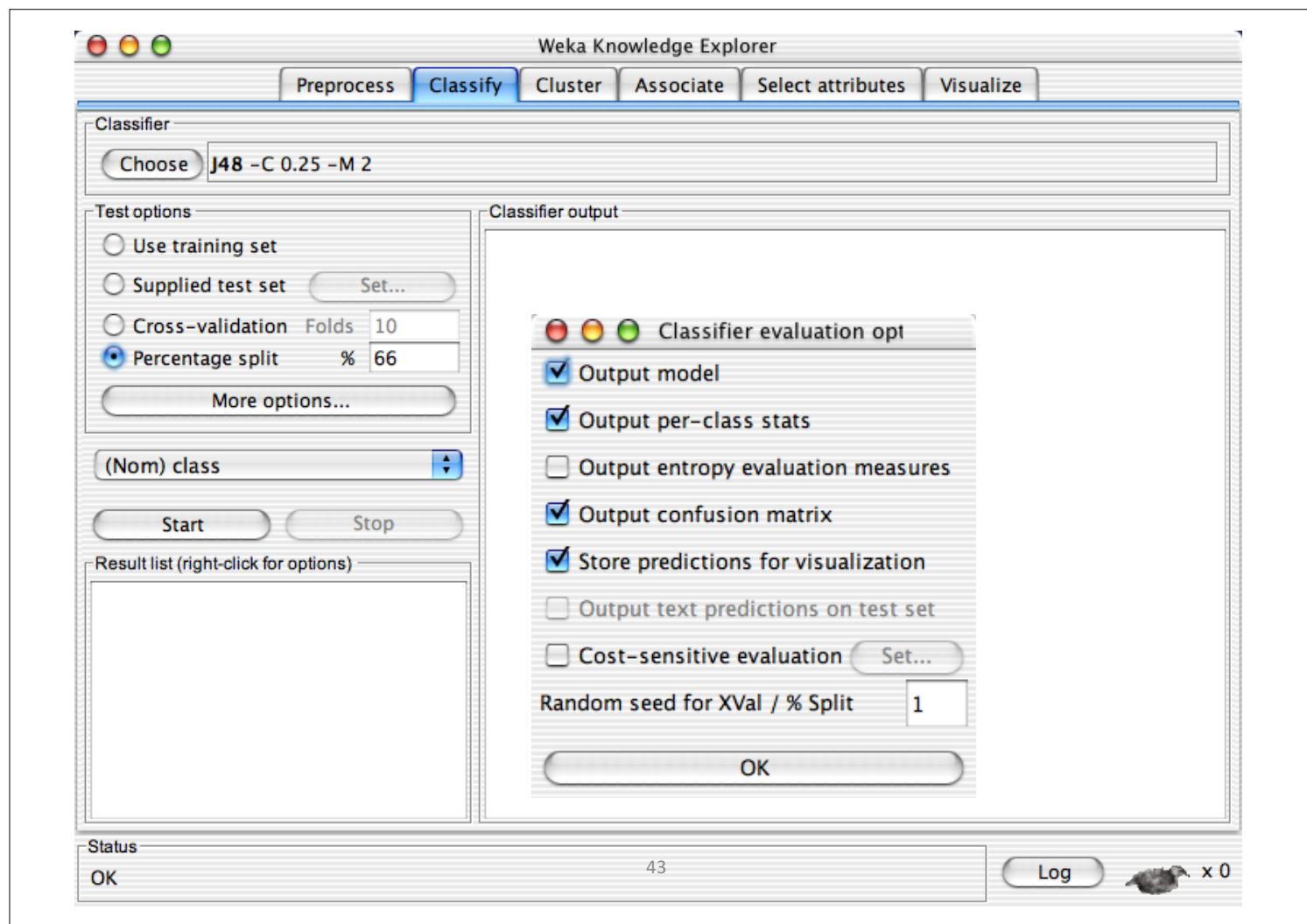


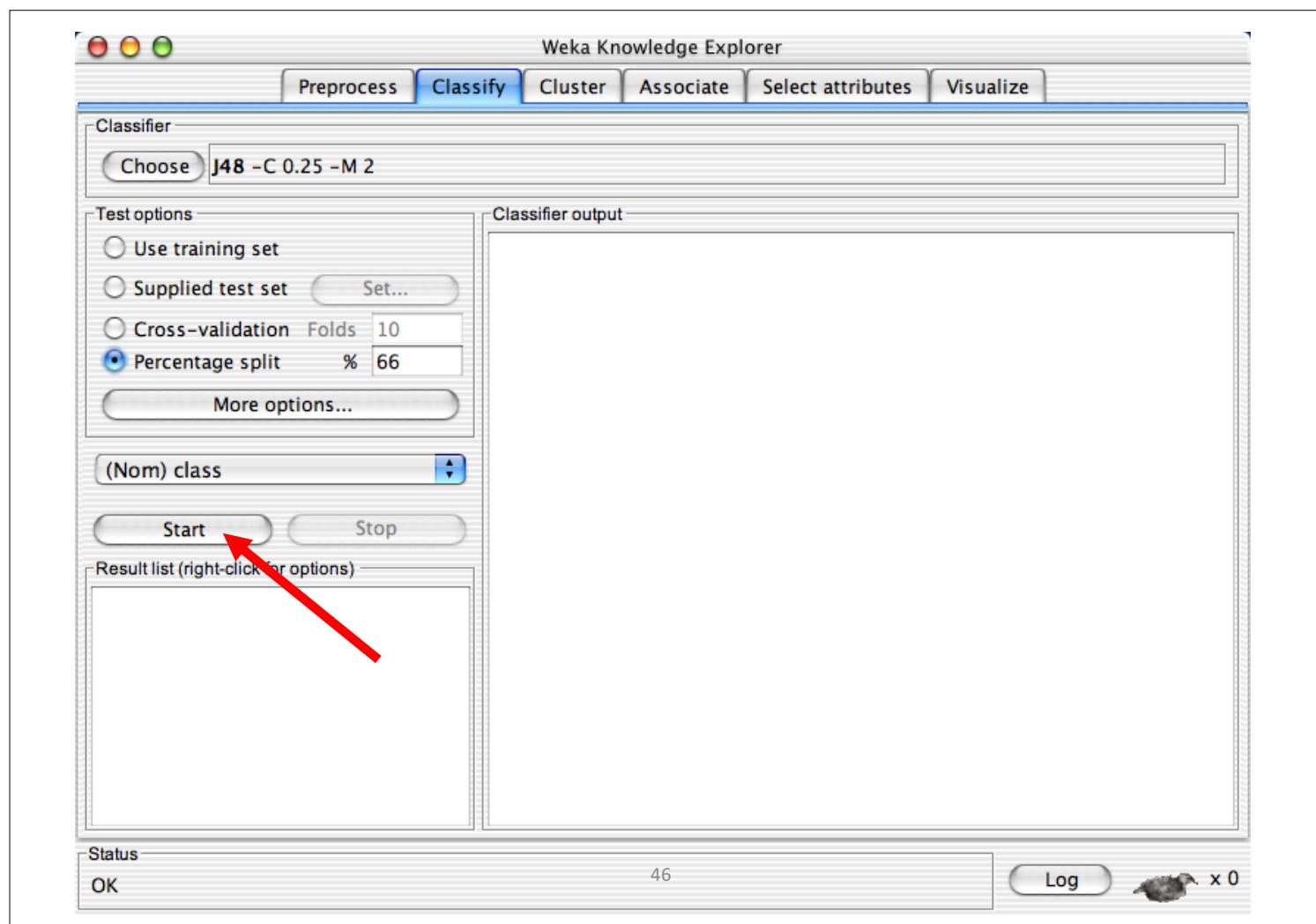
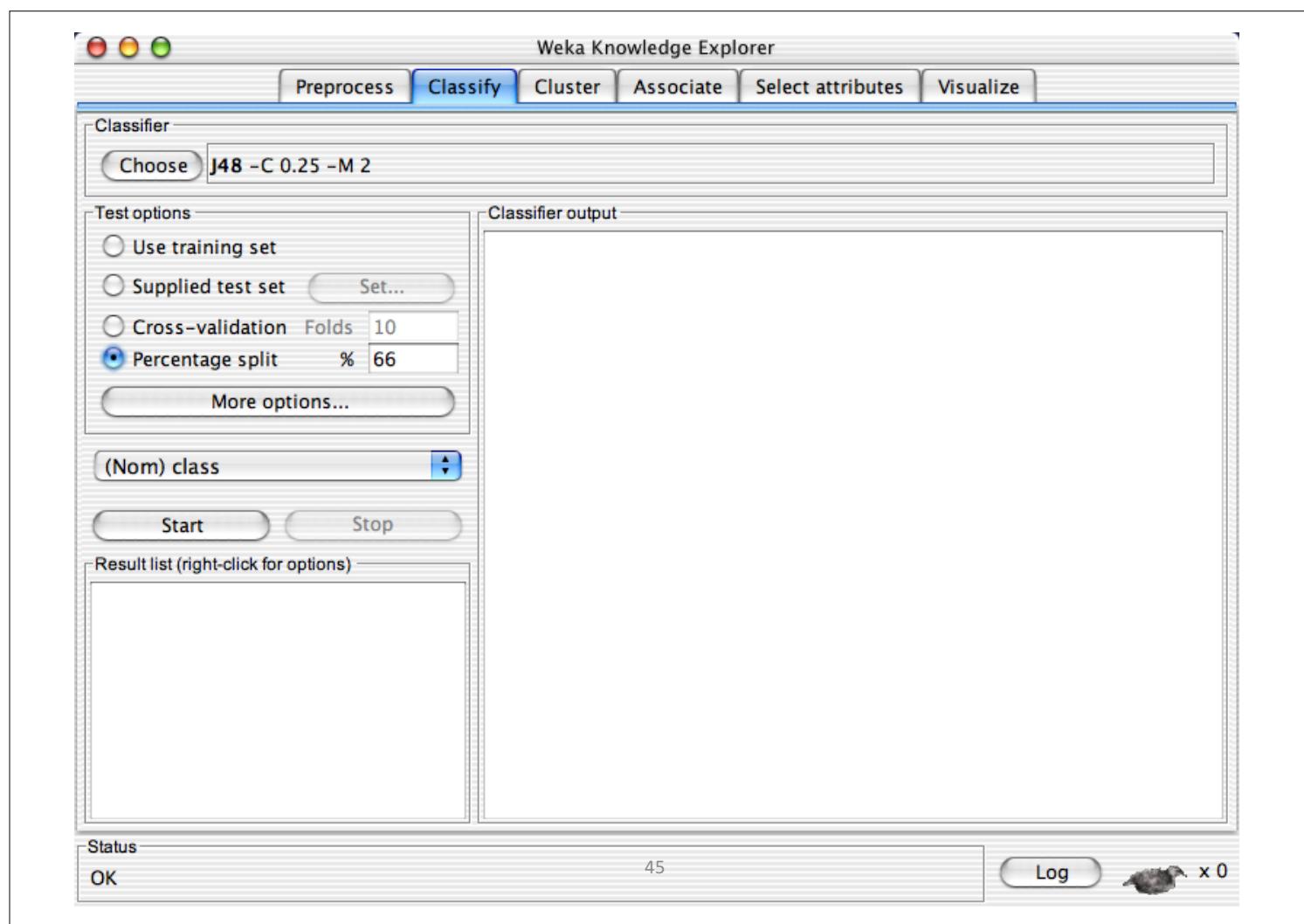












Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 - C 0.25 - M 2

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) class

Start Stop

Result list (right-click for options)

11:49:05 - trees.j48.J48

Classifier output

```
== Run information ==

Scheme: weka.classifiers.trees.J48 -C 0.25 -M 2
Relation: iris
Instances: 150
Attributes: 5
sepallength
sepalwidth
petallength
petalwidth
class
Test mode: split 66% train, remainder test

== Classifier model (full training set) ==
J48 pruned tree
-----
petalwidth <= 0.6: Iris-setosa (50.0)
petalwidth > 0.6
| petalwidth <= 1.7
| | petallength <= 4.9: Iris-versicolor (48.0/1.0)
| | petallength > 4.9
| | | petalwidth <= 1.5: Iris-virginica (3.0)
| | | petalwidth > 1.5: Iris-versicolor (3.0/1.0)
| petalwidth > 1.7: Iris-virginica (46.0/1.0)

Number of Leaves : 5
```

Status

OK 47 Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 - C 0.25 - M 2

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) class

Start Stop

Result list (right-click for options)

11:49:05 - trees.j48.J48

Classifier output

```
== Run information ==

Scheme: weka.classifiers.trees.J48 -C 0.25 -M 2
Relation: iris
Instances: 150
Attributes: 5
sepallength
sepalwidth
petallength
petalwidth
class
Test mode: split 66% train, remainder test

== Classifier model (full training set) ==
J48 pruned tree
-----
petalwidth <= 0.6: Iris-setosa (50.0)
petalwidth > 0.6
| petalwidth <= 1.7
| | petallength <= 4.9: Iris-versicolor (48.0/1.0)
| | petallength > 4.9
| | | petalwidth <= 1.5: Iris-virginica (3.0)
| | | petalwidth > 1.5: Iris-versicolor (3.0/1.0)
| petalwidth > 1.7: Iris-virginica (46.0/1.0)

Number of Leaves : 5
```

Status

OK 48 Log x 0



Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 - C 0.25 - M 2

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) class

Start Stop

Result list (right-click for options)

11:49:05 - trees.j48.J48

Classifier output

```
Time taken to build model: 0.24 seconds
===
Evaluation on test split ===
===
Summary ===

Correctly Classified Instances      49          96.0784 %
Incorrectly Classified Instances   2           3.9216 %
Kappa statistic                   0.9408
Mean absolute error               0.0396
Root mean squared error          0.1579
Relative absolute error          8.8979 %
Root relative squared error     33.4091 %
Total Number of Instances        51

===
Detailed Accuracy By Class ===

TP Rate    FP Rate    Precision    Recall    F-Measure    Class
1          0          1            1          1          Iris-setosa
1          0.063      0.905       1          0.95       Iris-versicolor
0.882     0          1            0.882     0.938      Iris-virginica

===
Confusion Matrix ===

a   b   c   <-- classified as
15  0   0   |   a = Iris-setosa
0  19  0   |   b = Iris-versicolor
0  2   15  |   c = Iris-virginica
```

Status

OK 49 Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 - C 0.25 - M 2

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) class

Start Stop

Result list (right-click for options)

11:49:05 - trees.j48.J48

Classifier output

```
Time taken to build model: 0.24 seconds
===
Evaluation on test split ===
===
Summary ===

Correctly Classified Instances      49          96.0784 %
Incorrectly Classified Instances   2           3.9216 %
Kappa statistic                   0.9408
Mean absolute error               0.0396
Root mean squared error          0.1579
Relative absolute error          8.8979 %
Root relative squared error     33.4091 %
Total Number of Instances        51

===
Detailed Accuracy By Class ===

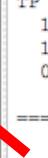
TP Rate    FP Rate    Precision    Recall    F-Measure    Class
1          0          1            1          1          Iris-setosa
1          0.063      0.905       1          0.95       Iris-versicolor
0.882     0          1            0.882     0.938      Iris-virginica

===
Confusion Matrix ===

a   b   c   <-- classified as
15  0   0   |   a = Iris-setosa
0  19  0   |   b = Iris-versicolor
0  2   15  |   c = Iris-virginica
```

Status

OK 50 Log x 0



Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 - C 0.25 - M 2

Test options

- Use training set
- Supplied test set [Set...](#)
- Cross-validation Folds 10
- Percentage split % 66

[More options...](#)

(Nom) class

[Start](#) [Stop](#)

Result list (right-click for options)

11:49:05 - trees.j48.J48

Classifier output

```
Time taken to build model: 0.24 seconds
===
Evaluation on test split ===
Summary ===

Correctly Classified Instances      49          96.0784 %
Incorrectly Classified Instances   2           3.9216 %
Kappa statistic                   0.9408
Mean absolute error               0.0396
Root mean squared error           0.1579
Relative absolute error           8.8979 %
Root relative squared error     33.4091 %
Total Number of Instances        51

===
Detailed Accuracy By Class ===
```

	Recall	F-Measure	Class
1	1	Iris-setosa	
1	0.95	Iris-versicolor	
0.882	0.938	Iris-virginica	

[View in main window](#)

[View in separate window](#)

[Save result buffer](#)

[Load model](#)

[Save model](#)

[Re-evaluate model on current test set](#)

[Visualize classifier errors](#)

Visualize tree

[Visualize margin curve](#)

[Visualize threshold curve 51](#)

[Visualize cost curve](#)

Status

OK

[Log](#) x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 - C 0.25 - M 2 Weka Classifier Tree Visualizer: 11:49:05 - trees.j48.J48 (iris)

Test options

- Use training set
- Supplied test set
- Cross-validation
- Percentage split

[More options...](#)

(Nom) class

[Start](#)

Result list (right-click for options)

11:49:05 - trees.j48.J48

Tree View

```

graph TD
    Root[petalwidth] --<= 0.6--> IrisSetosa[Iris-setosa (50.0)]
    Root --> Petalwidth1[petalwidth]
    Petalwidth1 --<= 1.7--> IrisVersicolor1[Iris-versicolor (48.0/1.0)]
    Petalwidth1 --> Petalwidth2[petalwidth]
    Petalwidth2 --<= 4.9--> IrisVersicolor2[Iris-virginica (3.0)]
    Petalwidth2 --> IrisVersicolor3[Iris-versicolor (3.0/1.0)]
  
```

96.0784 %
3.9216 %

ass
is-setosa
is-versicolor
is-virginica

15 0 0 | a = Iris-setosa
0 19 0 | b = Iris-versicolor
0 2 15 | c = Iris-virginica

Status

OK

[Log](#) x 0

Weka Knowledge Explorer

Classify (selected)

Classifier

Choose: J48 - C 0.25 - M 2

Test options

- Use training set
- Supplied test set [Set...](#)
- Cross-validation Folds 10
- Percentage split % 66

[More options...](#)

(Nom) class

Start **Stop**

Result list (right-click for options)

11:49:05 - trees.j48.J48

- [View in main window](#)
- [View in separate window](#)
- [Save result buffer](#)
- [Load model](#)
- [Save model](#)
- [Re-evaluate model on current test set](#)
- [Visualize classifier errors](#) (selected)
- [Visualize tree](#)
- [Visualize margin curve](#)
- [Visualize threshold curve](#) 53
- [Visualize cost curve](#)

Classifier output

```
Time taken to build model: 0.24 seconds
===
Evaluation on test split ===
Summary ===

Correctly Classified Instances      49          96.0784 %
Incorrectly Classified Instances   2           3.9216 %
Kappa statistic                   0.9408
Mean absolute error               0.0396
Root mean squared error          0.1579
Relative absolute error          8.8979 %
Root relative squared error     33.4091 %
Total Number of Instances        51

Detailed Accuracy By Class ===

Recall    F-Measure    Class
1         1           Iris-setosa
1         0.95        Iris-versicolor
0.882    0.938       Iris-virginica
```

Status

OK

Log x 0

Weka Knowledge Explorer

Classify (selected)

Classifier

Choose: J48 - C 0.25 - M 2

Test options

- Use training set
- Supplied test set
- Cross-validation
- Percentage split

X: petallength (Num) Y: petalwidth (Num)
Colour: class (Nom) Select Instance

[Reset](#) [Clear](#) [Save](#) [Jitter](#)

(Nom) class

Start

Result list (right-click for options)

11:49:05 - trees.j48.J48

Weka Classifier Visualize: 11:49:05 - trees.j48.J48 (iris)

Plot: iris_predicted

Class colour

Iris-setosa	Iris-versicolor	Iris-virginica
-------------	-----------------	----------------

```
v t v u - Iris-versicolor
0 2 15 | c = Iris-virginica
```

Status

OK

Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 - C 0.25 - M 2

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) class

Start Stop

Result list (right-click for options)

11:49:05 - trees.j48.J48

Classifier output

```
Time taken to build model: 0.24 seconds
===
Evaluation on test split ===
===
Summary ===

Correctly Classified Instances      49          96.0784 %
Incorrectly Classified Instances   2           3.9216 %
Kappa statistic                   0.9408
Mean absolute error               0.0396
Root mean squared error          0.1579
Relative absolute error          8.8979 %
Root relative squared error     33.4091 %
Total Number of Instances        51

===
Detailed Accuracy By Class ===

TP Rate    FP Rate    Precision    Recall    F-Measure    Class
1          0          1            1          1          Iris-setosa
1          0.063      0.905       1          0.95       Iris-versicolor
0.882      0          1            0.882     0.938       Iris-virginica

===
Confusion Matrix ===

a   b   c   <-- classified as
15  0   0   |  a = Iris-setosa
0   19  0   |  b = Iris-versicolor
0   0   215 |  c = Iris-virginica
```

Status

OK 55 Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 - C 0.25 - M 2

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) class

Start Stop

Result list (right-click for options)

11:49:05 - trees.j48.J48

Classifier output

```
Time taken to build model: 0.24 seconds
===
Evaluation on test split ===
===
Summary ===

Correctly Classified Instances      49          96.0784 %
Incorrectly Classified Instances   2           3.9216 %
Kappa statistic                   0.9408
Mean absolute error               0.0396
Root mean squared error          0.1579
Relative absolute error          8.8979 %
Root relative squared error     33.4091 %
Total Number of Instances        51

===
Detailed Accuracy By Class ===

TP Rate    FP Rate    Precision    Recall    F-Measure    Class
1          0          1            1          1          Iris-setosa
1          0.063      0.905       1          0.95       Iris-versicolor
0.882      0          1            0.882     0.938       Iris-virginica

===
Confusion Matrix ===

a   b   c   <-- classified as
15  0   0   |  a = Iris-setosa
0   19  0   |  b = Iris-versicolor
0   0   215 |  c = Iris-virginica
```

Status

OK 56 Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

- weka
- classifiers
 - bayes
 - AODE
 - BayesNetK2
 - BayesNetB
 - NaiveBayes**
 - NaiveBayesMultinomial
 - NaiveBayesSimple
 - NaiveBayesUpdateable
 - functions
 - lazy
 - meta
 - misc
 - trees
 - rules

Classifier output

```
== Evaluation on test split ==
== Summary ==

Correctly Classified Instances      50          98.0392 %
Incorrectly Classified Instances    1           1.9608 %
Kappa statistic                      0.9704
Mean absolute error                  0.0239
Root mean squared error              0.1101 %
Relative absolute error              5.3594 %
Root relative squared error         23.2952 %
Total Number of Instances            51

== Detailed Accuracy By Class ==

      TP Rate   FP Rate   Precision   Recall   F-Measure   Class
      1          0          1          1          1          Iris-setosa
      1          0.031     0.95       1          0.974     Iris-versicolor
      0.941     0          1          0.941     0.97       Iris-virginica

== Confusion Matrix ==

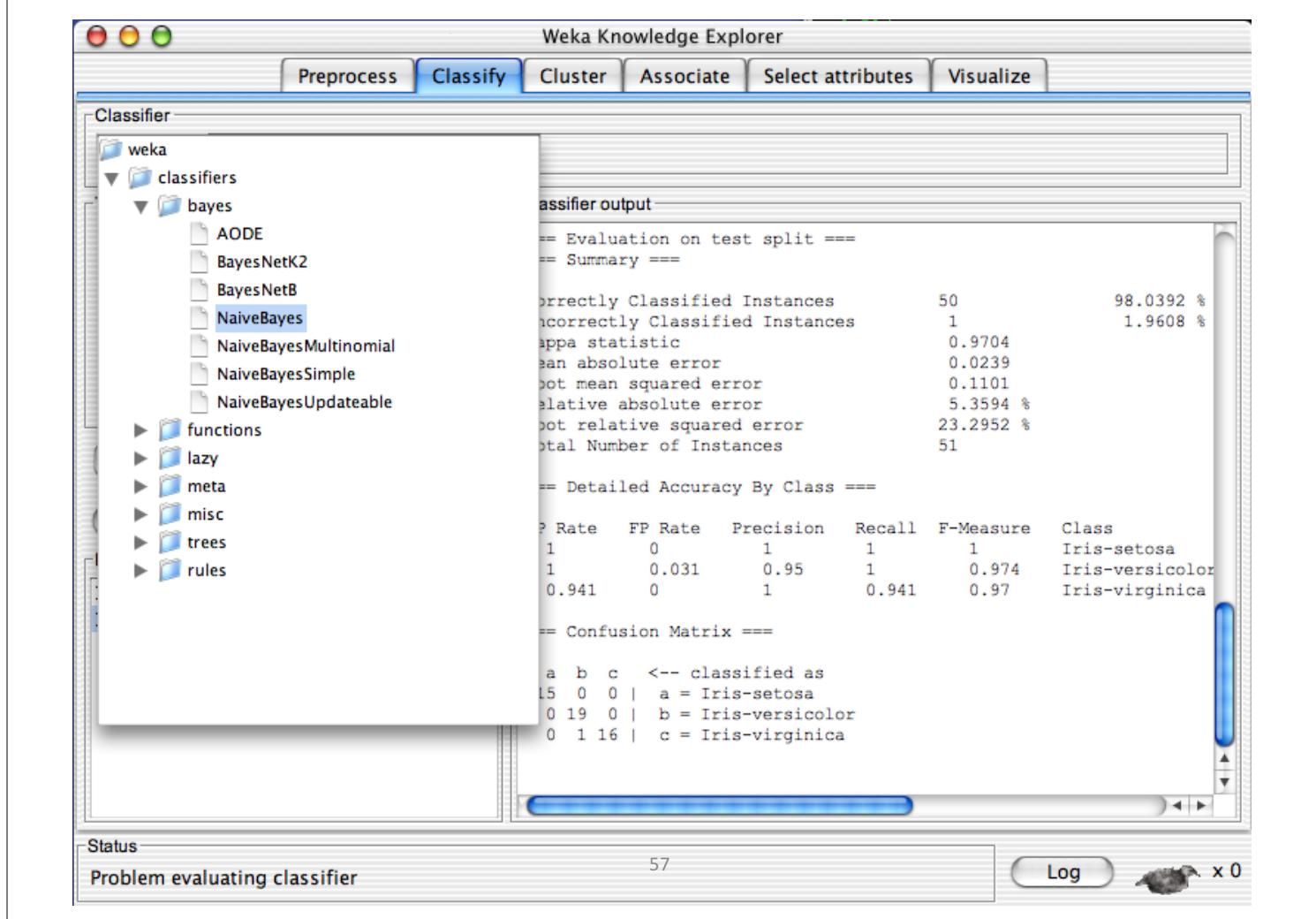
      a   b   c   <-- classified as
15   0   0   |   a = Iris-setosa
 0  19   0   |   b = Iris-versicolor
 0   1  16   |   c = Iris-virginica
```

Status

Problem evaluating classifier

57

Log x 0



Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **NaiveBayes**

Test options

Use training set

Supplied test set Set...

Cross-validation Folds 10

Percentage split % 66

More options...

(Nom) class

Start Stop

Result list (right-click for options)

- 11:49:05 - trees.j48.J48
- 14:34:28 - functions.neural.NeuralNetwork**

Classifier output

```
== Evaluation on test split ==
== Summary ==

Correctly Classified Instances      50          98.0392 %
Incorrectly Classified Instances    1           1.9608 %
Kappa statistic                      0.9704
Mean absolute error                  0.0239
Root mean squared error              0.1101 %
Relative absolute error              5.3594 %
Root relative squared error         23.2952 %
Total Number of Instances            51

== Detailed Accuracy By Class ==

      TP Rate   FP Rate   Precision   Recall   F-Measure   Class
      1          0          1          1          1          Iris-setosa
      1          0.031     0.95       1          0.974     Iris-versicolor
      0.941     0          1          0.941     0.97       Iris-virginica

== Confusion Matrix ==

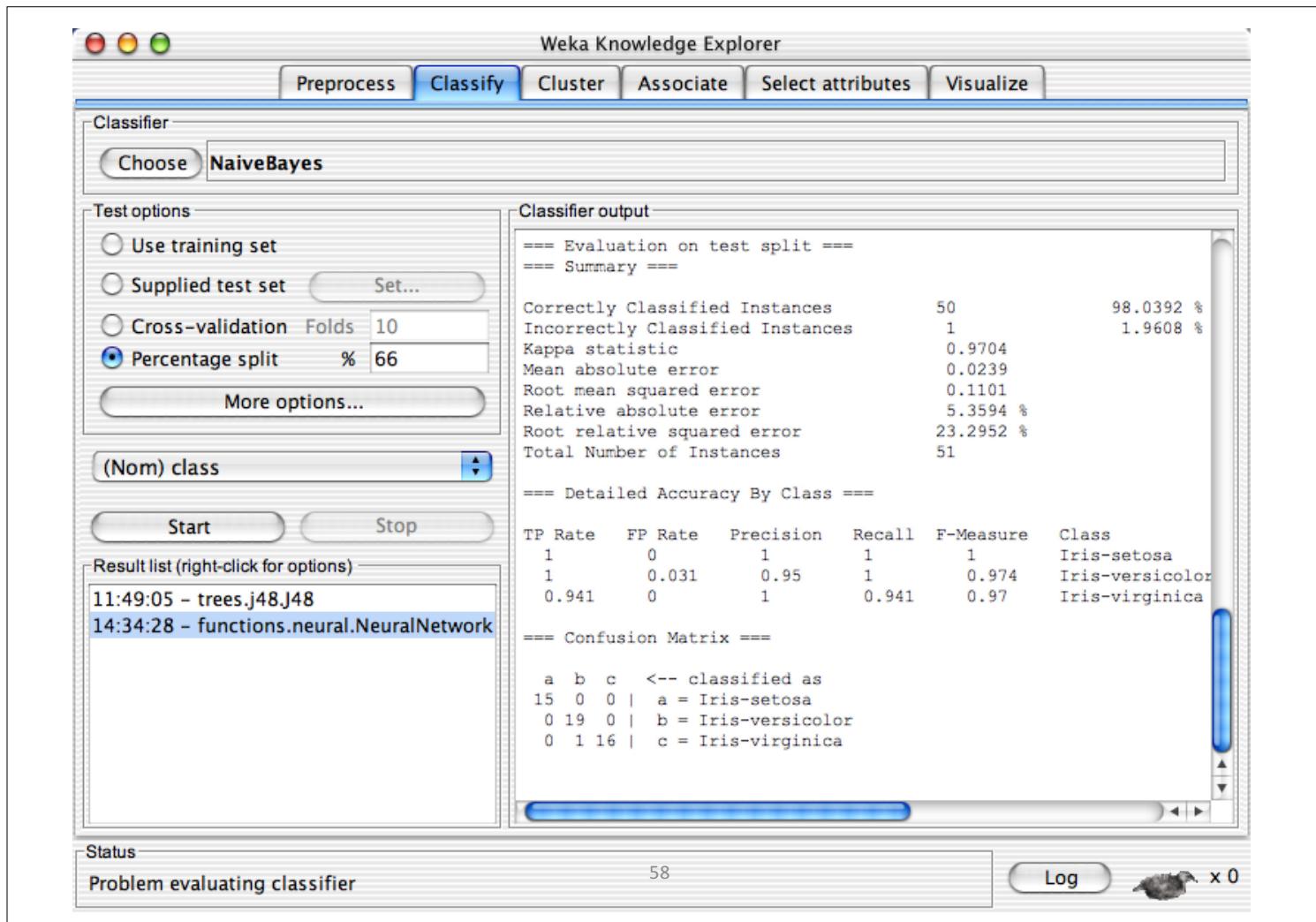
      a   b   c   <-- classified as
15   0   0   |   a = Iris-setosa
 0  19   0   |   b = Iris-versicolor
 0   1  16   |   c = Iris-virginica
```

Status

Problem evaluating classifier

58

Log x 0



Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **NaiveBayes**

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) class

Start Stop

Result list (right-click for options)

11:49:05 - trees.j48.J48
14:34:28 - functions.neural.NeuralNetwork

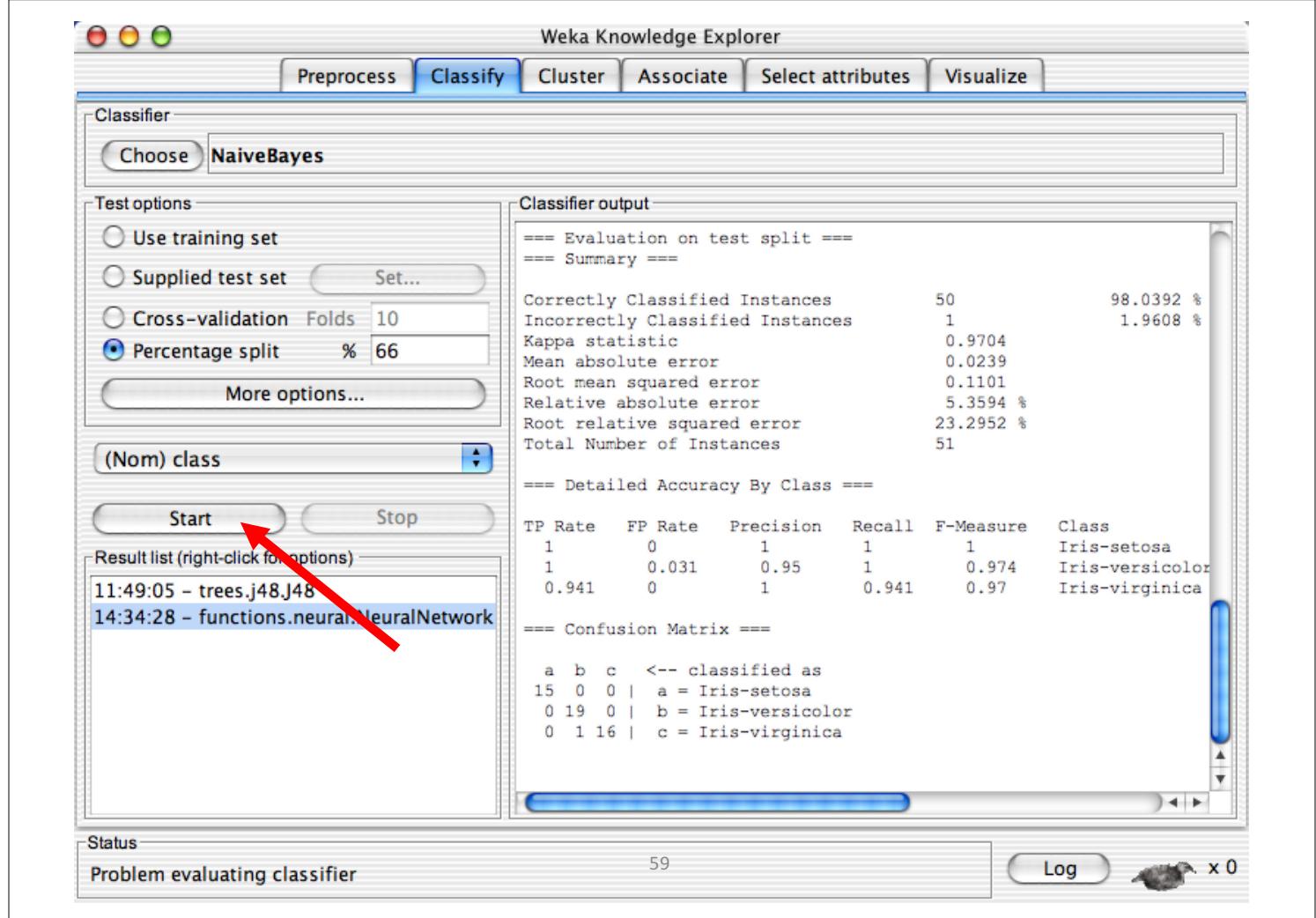
Classifier output

==== Evaluation on test split ====
==== Summary ====
Correctly Classified Instances 50 98.0392 %
Incorrectly Classified Instances 1 1.9608 %
Kappa statistic 0.9704
Mean absolute error 0.0239
Root mean squared error 0.1101
Relative absolute error 5.3594 %
Root relative squared error 23.2952 %
Total Number of Instances 51

==== Detailed Accuracy By Class ====
TP Rate FP Rate Precision Recall F-Measure Class
1 0 1 1 1 Iris-setosa
1 0.031 0.95 1 0.974 Iris-versicolor
0.941 0 1 0.941 0.97 Iris-virginica

==== Confusion Matrix ====
a b c <-- classified as
15 0 0 | a = Iris-setosa
0 19 0 | b = Iris-versicolor
0 1 16 | c = Iris-virginica

Status
Problem evaluating classifier 59 Log x 0



Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **NaiveBayes**

Test options

Use training set
 Supplied test set Set...
 Cross-validation Folds 10
 Percentage split % 66
More options...

(Nom) class

Start Stop

Result list (right-click for options)

11:49:05 - trees.j48.J48
14:34:28 - functions.neural.NeuralNetwork
14:48:05 - bayes.NaiveBayes

Classifier output

==== Evaluation on test split ====
==== Summary ====
Correctly Classified Instances 48 94.1176 %
Incorrectly Classified Instances 3 5.8824 %
Kappa statistic 0.9113
Mean absolute error 0.0447
Root mean squared error 0.1722
Relative absolute error 10.0365 %
Root relative squared error 36.4196 %
Total Number of Instances 51

==== Detailed Accuracy By Class ====
TP Rate FP Rate Precision Recall F-Measure Class
1 0 1 1 1 Iris-setosa
0.947 0.063 0.9 0.947 0.923 Iris-versicolor
0.882 0.029 0.938 0.882 0.909 Iris-virginica

==== Confusion Matrix ====
a b c <-- classified as
15 0 0 | a = Iris-setosa
0 18 1 | b = Iris-versicolor
0 2 15 | c = Iris-virginica

Status
OK 60 Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **NaiveBayes**

Test options

- Use training set
- Supplied test set Set...
- Cross-validation Folds 10
- Percentage split % 66

(Nom) class

Start Stop

Result list (right-click for options)

- 11:49:05 - trees.j48.J48
- 14:34:28 - functions.neural.NeuralNetwork
- 14:48:05 - bayes.NaiveBayes**

Classifier output

```
==== Evaluation on test split ====
==== Summary ===

Correctly Classified Instances      48          94.1176 %
Incorrectly Classified Instances   3           5.8824 %
Kappa statistic                   0.9113
Mean absolute error               0.0447
Root mean squared error           0.1722
Relative absolute error           10.0365 %
Root relative squared error      36.4196 %
Total Number of Instances         51

==== Detailed Accuracy By Class ===

TP Rate   FP Rate   Precision   Recall   F-Measure   Class
1        0         1           1         1           Iris-setosa
0.947    0.063     0.9         0.947    0.923      Iris-versicolor
0.882    0.029     0.938       0.882    0.909      Iris-virginica
```

```
==== Confusion Matrix ===

a b c    <-- classified as
15 0 0 | a = Iris-setosa
 0 18 1 | b = Iris-versicolor
 0 2 15 | c = Iris-virginica
```

Status

OK 61 Log x 0

Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose **NaiveBayes**

Test options

- Use training set
- Supplied test set Set...
- Cross-validation Folds 10
- Percentage split % 66

(Nom) class

Start View in main window View in separate window Save result buffer

Result list (right-click for options)

- 11:49:05 - trees.j48.J48
- 14:34:28 - functions.NeuralNetwork
- 14:48:05 - bayes.NaiveBayes**

Classifier output

```
==== Evaluation on test split ====
==== Summary ===

Correctly Classified Instances      48          94.1176 %
Incorrectly Classified Instances   3           5.8824 %
Kappa statistic                   0.9113
Mean absolute error               0.0447
Root mean squared error           0.1722
Relative absolute error           10.0365 %
Root relative squared error      36.4196 %
Total Number of Instances         51

==== Detailed Accuracy By Class ===

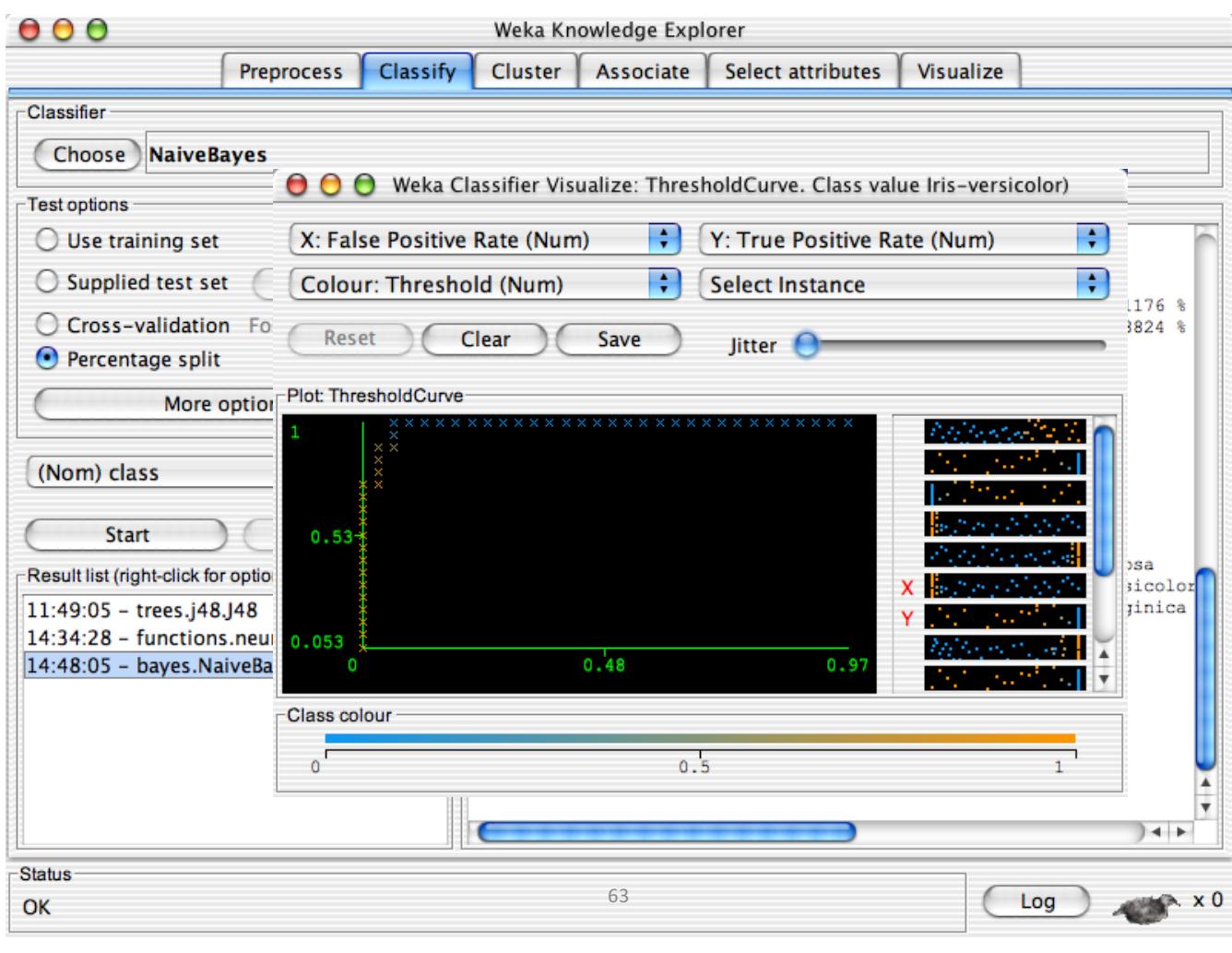
Precision   Recall   F-Measure   Class
1          1         1           Iris-setosa
0.9        0.947    0.923      Iris-versicolor
0.938      0.882    0.909      Iris-virginica
```

```
==== Confusion Matrix ===

a b c    <-- classified as
15 0 0 | a = Iris-setosa
 0 18 1 | b = Iris-versicolor
 0 2 15 | c = Iris-virginica
```

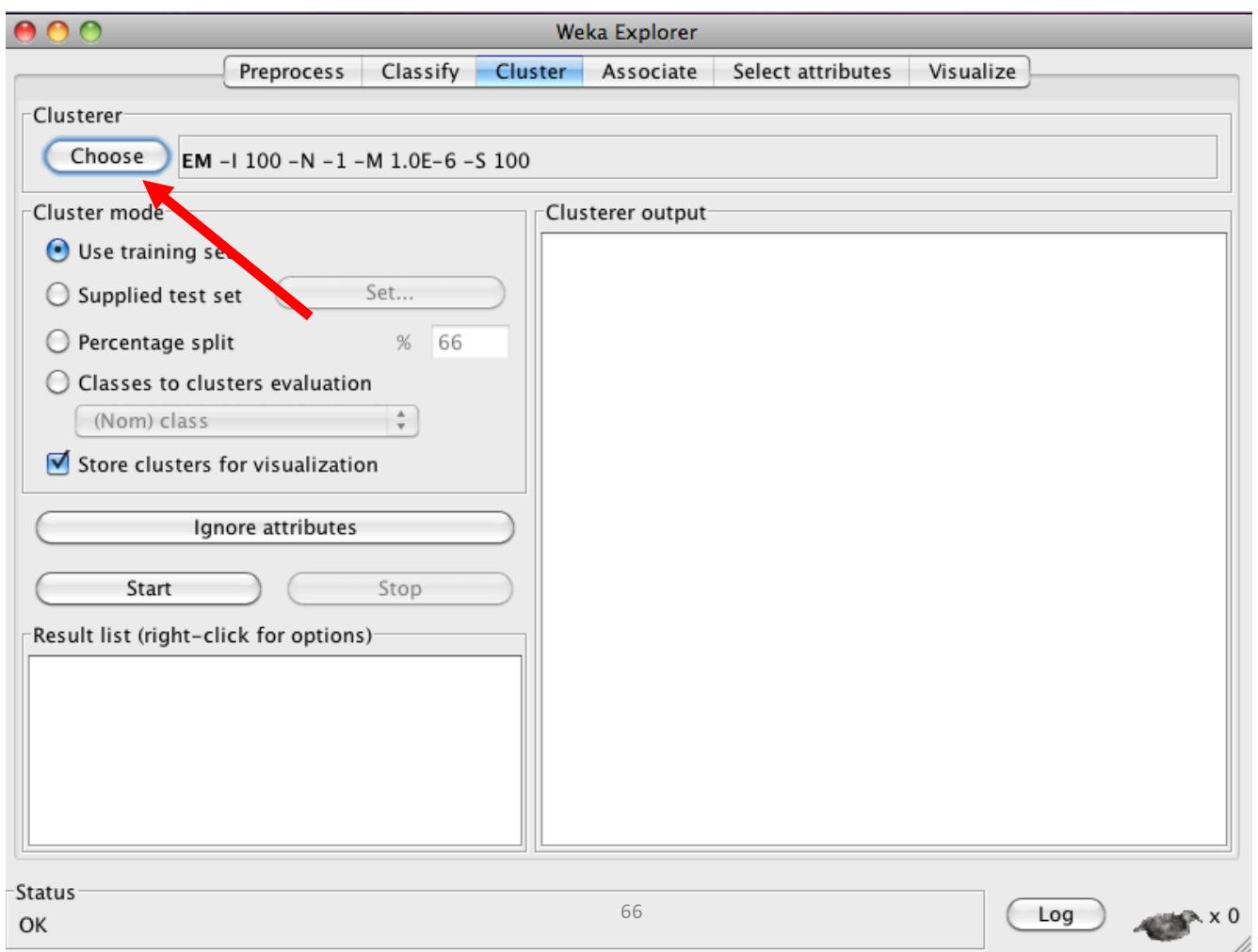
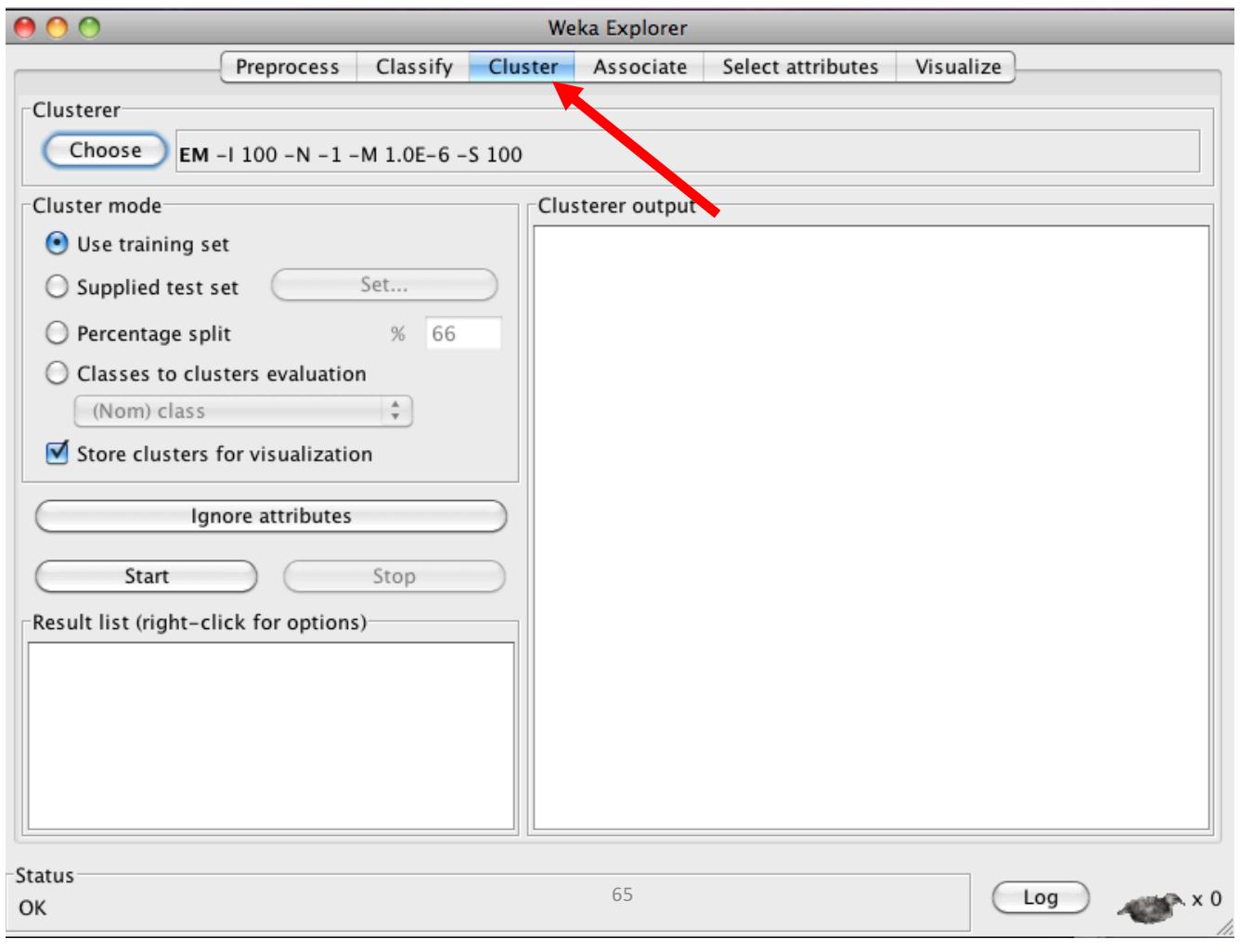
Status

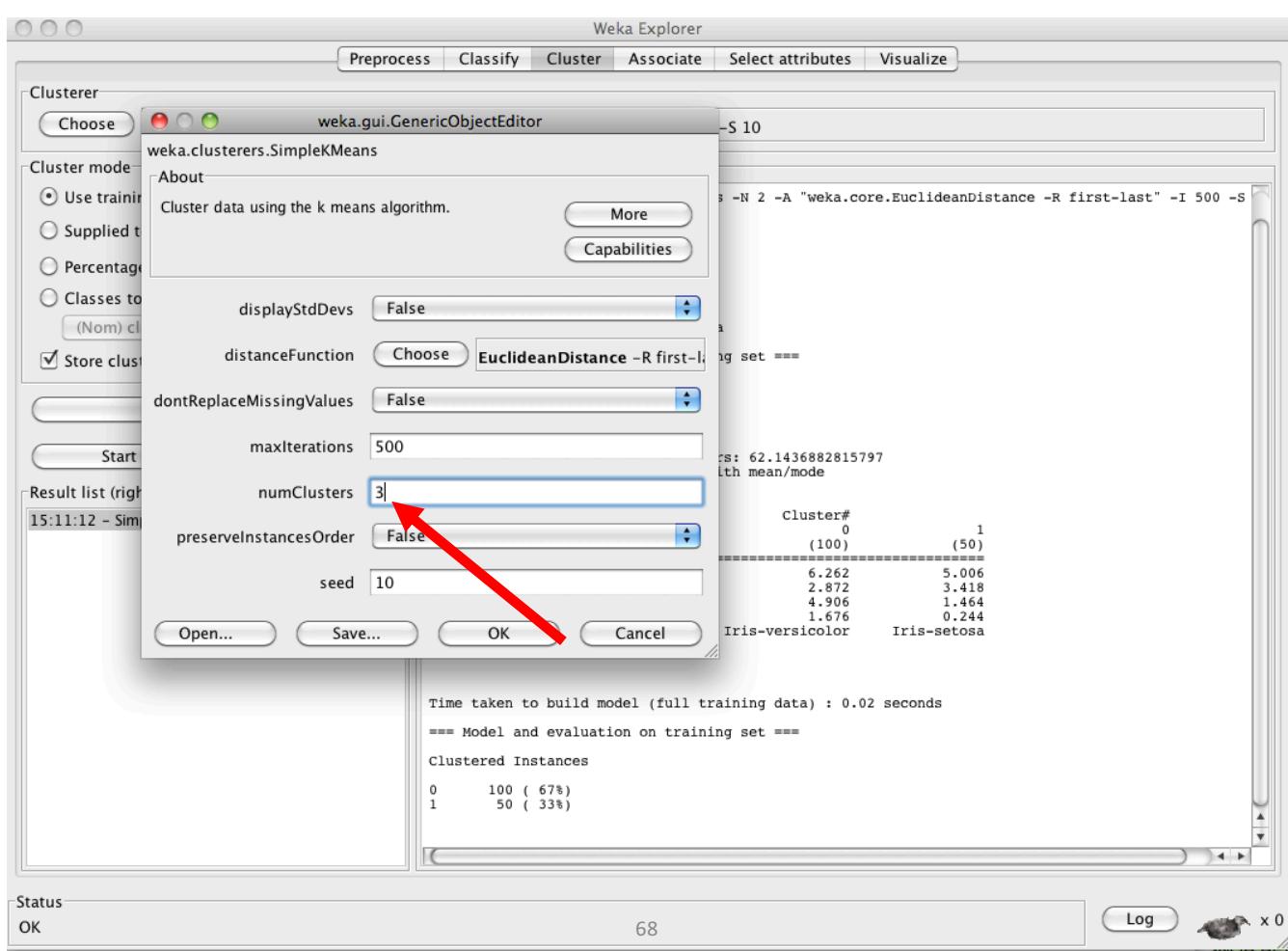
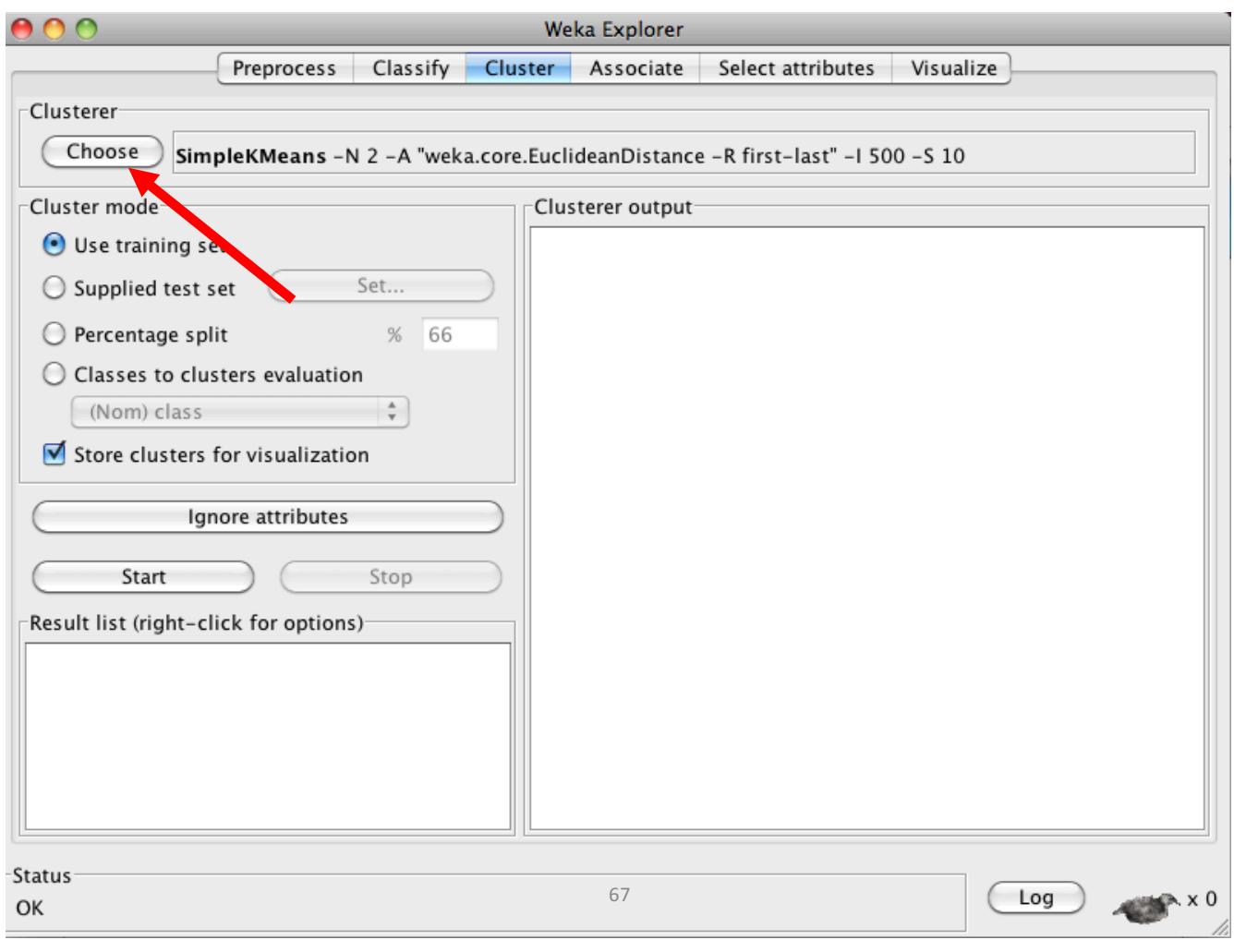
OK 62 Log x 0



Clustering

- WEKA contains many clustering implementations:
 - Works with both discrete and numerical data
- Example of K-means





Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Clusterer

Choose SimpleKMeans -N 3 -A "weka.core.EuclideanDistance -R first-last" -I 500 -S 10

Cluster mode

- Use training set
- Supplied test set Set...
- Percentage split % 66
- Classes to clusters evaluation (Nom) class
- Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

- 15:11:12 - SimpleKMeans
- 15:12:39 - SimpleKMeans**

Clusterer output

```

Relation: iris
Instances: 150
Attributes: 5
sepallength
sepalwidth
petallength
petalwidth
class

Test mode: evaluate on training data
*** Model and evaluation on training set ***
kMeans
=====

Number of iterations: 3
Within cluster sum of squared errors: 7.817456892309574
Missing values globally replaced with mean/mode

Cluster centroids:
```

Attribute	Full Data (150)	Cluster# 0 (50)	1 (50)	2 (50)
sepallength	5.8433	5.936	5.006	6.588
sepalwidth	3.054	2.77	3.418	2.974
petallength	3.7587	4.26	1.464	5.552
petalwidth	1.1987	1.326	0.244	2.026
class	Iris-setosa Iris-versicolor	Iris-setosa	Iris-virginica	

```

Time taken to build model (full training data) : 0 seconds
*** Model and evaluation on training set ***
Clustered Instances
0      50 ( 33%)
1      50 ( 33%)
2      50 ( 33%)
```

Status OK 69 Log x 0

Weka Clusterer Visualize: 15:12:39 - SimpleKMeans (iris)

X: Instance_number (Num) Y: sepallength (Num)

Colour: Cluster (Nom) Select Instance

Reset Clear Open Save Jitter

Plot:iris_clustered

Class colour

cluster0 cluster1 cluster2

Clusterer output

```

Relation: iris
Instances: 150
Attributes: 5
sepallength
sepalwidth
petallength
petalwidth
class

Test mode: evaluate on training data
*** Model and evaluation on training set ***
kMeans
=====

Number of iterations: 3
Within cluster sum of squared errors: 7.817456892309574
Missing values globally replaced with mean/mode

Cluster centroids:
```

Attribute	Full Data (150)	Cluster# 0 (50)	1 (50)	2 (50)
sepallength	5.8433	5.936	5.006	6.588
sepalwidth	3.054	2.77	3.418	2.974
petallength	3.7587	4.26	1.464	5.552
petalwidth	1.1987	1.326	0.244	2.026
class	Iris-setosa Iris-versicolor	Iris-setosa	Iris-virginica	

```

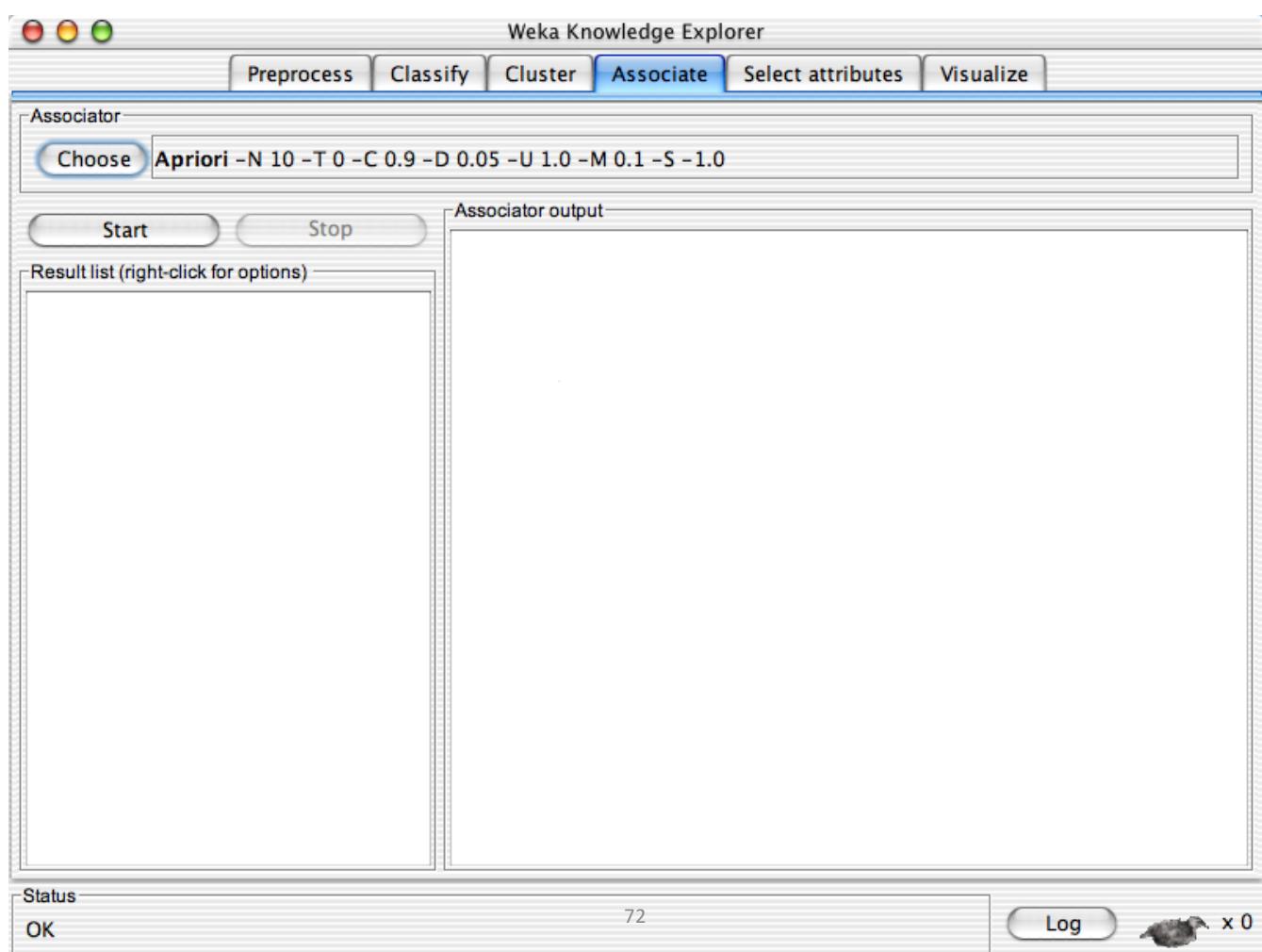
Time taken to build model (full training data) : 0 seconds
*** Model and evaluation on training set ***
Clustered Instances
0      50 ( 33%)
1      50 ( 33%)
2      50 ( 33%)
```

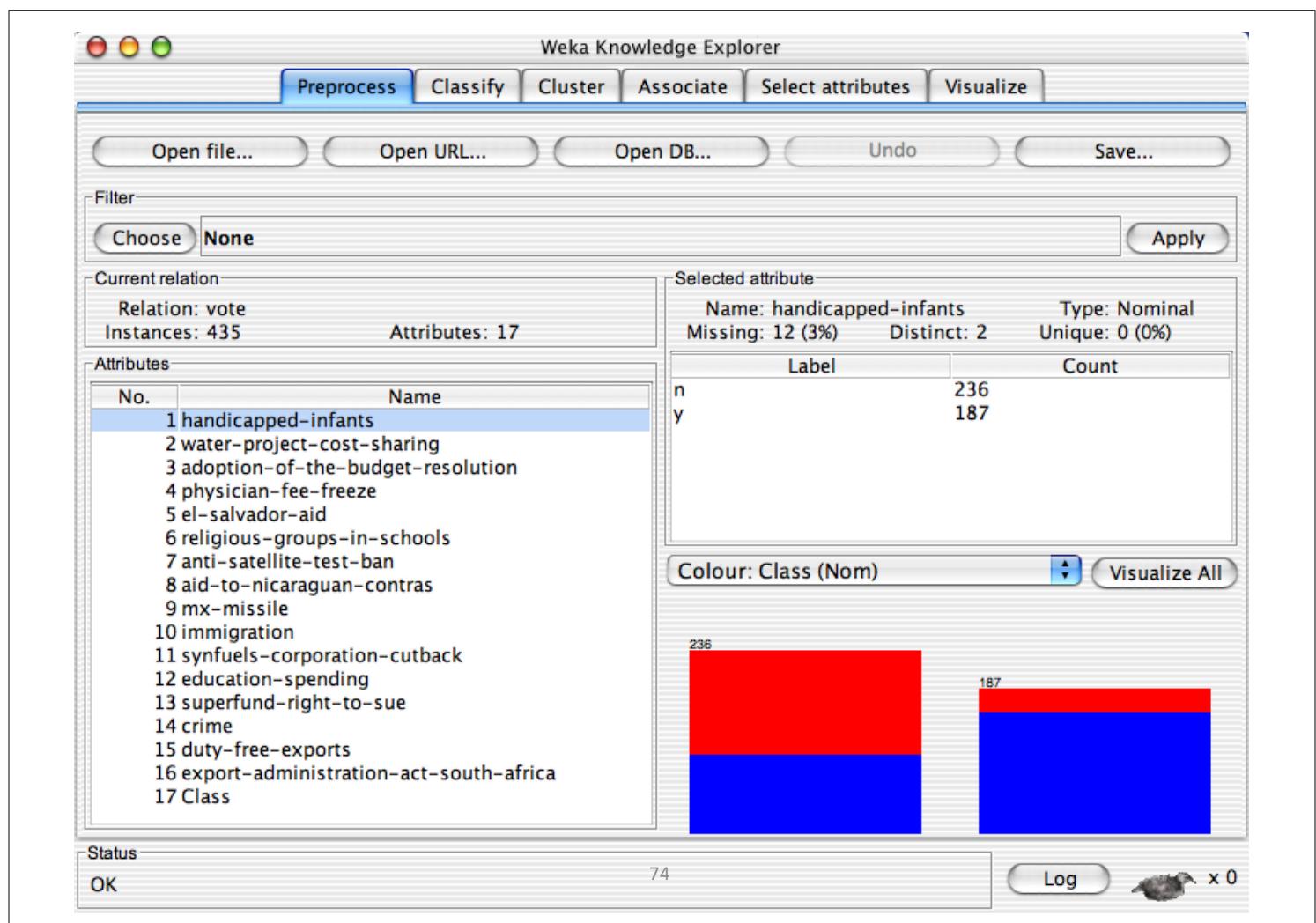
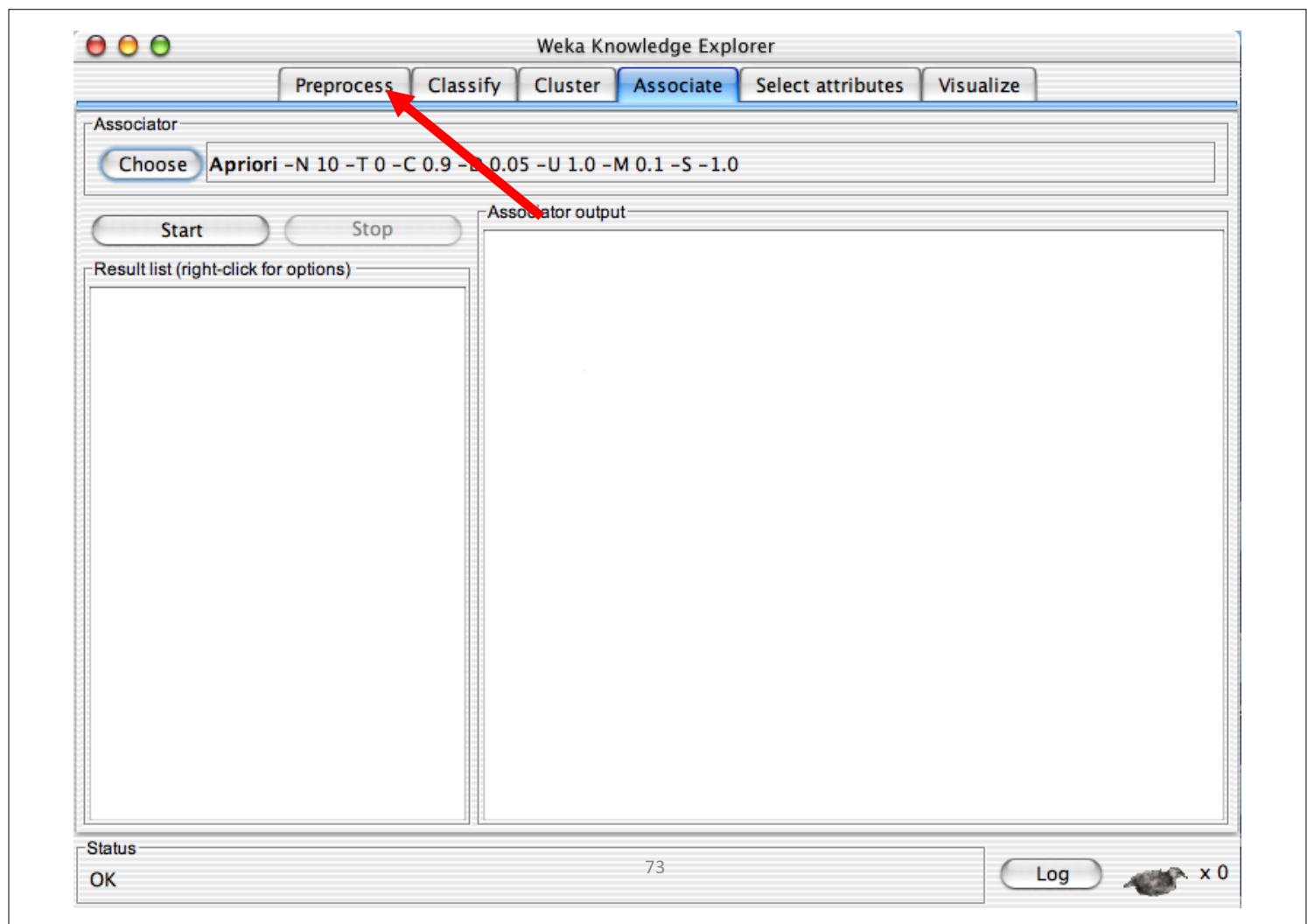
Status OK 70 Log x 0

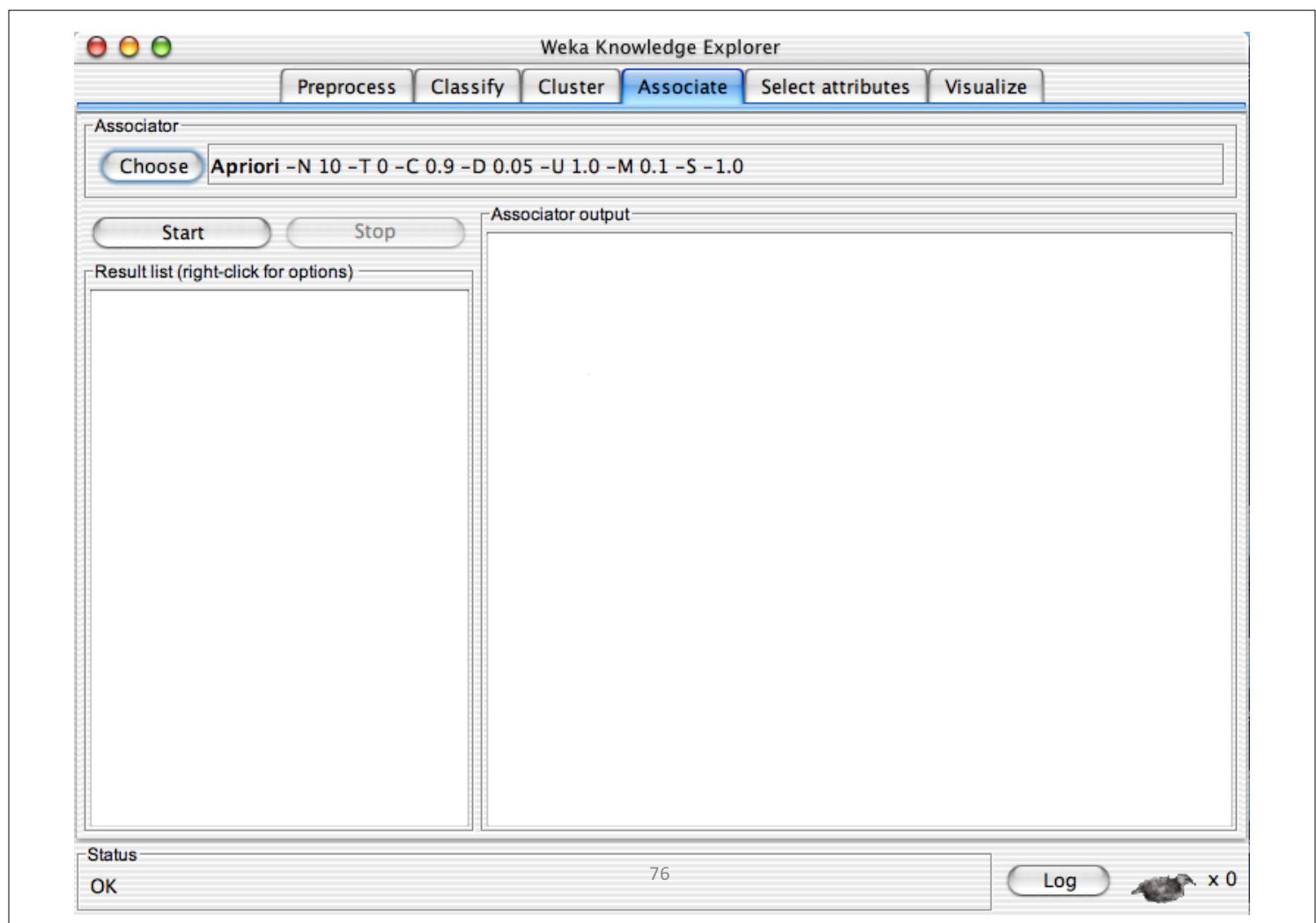
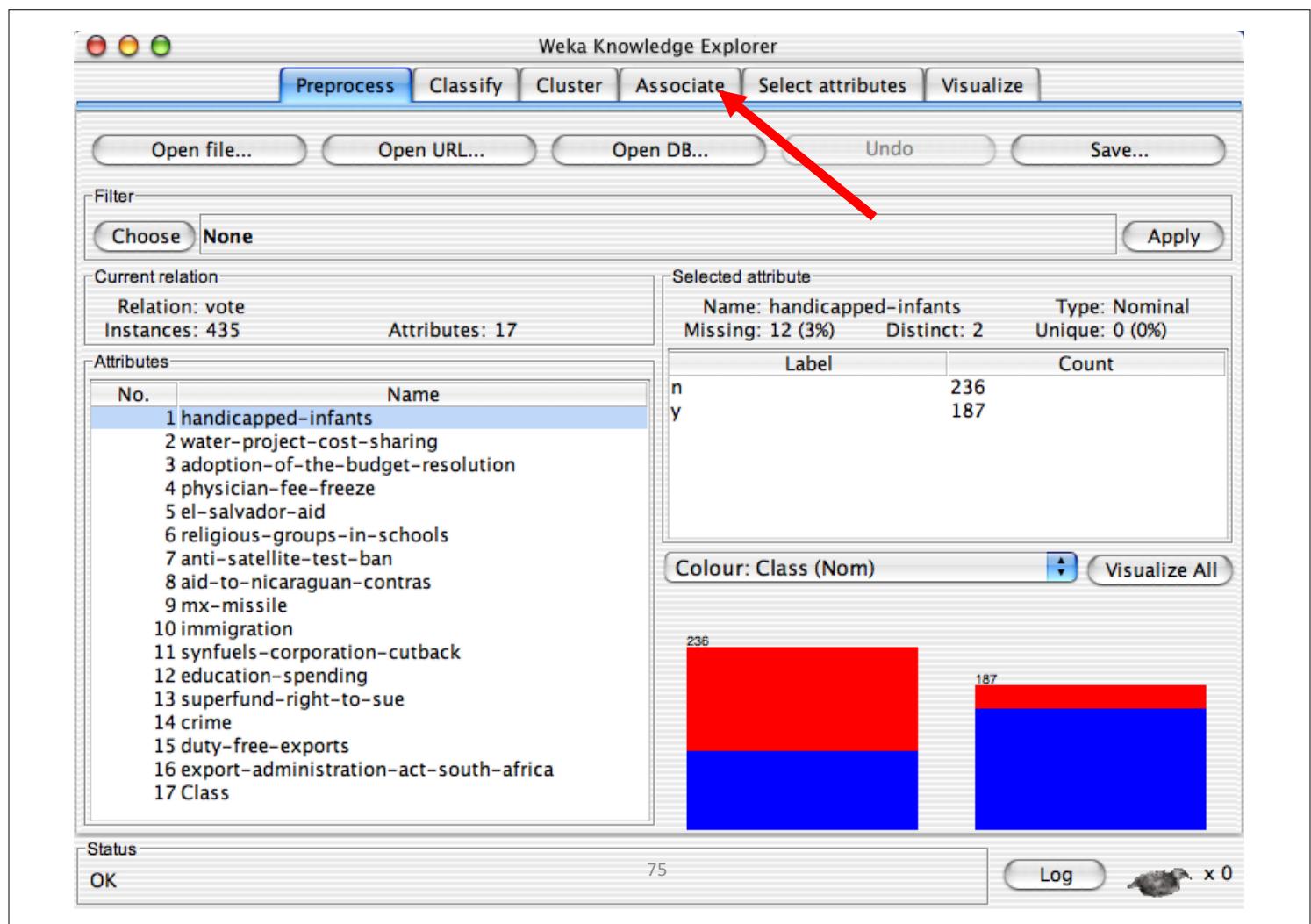
Finding Associations

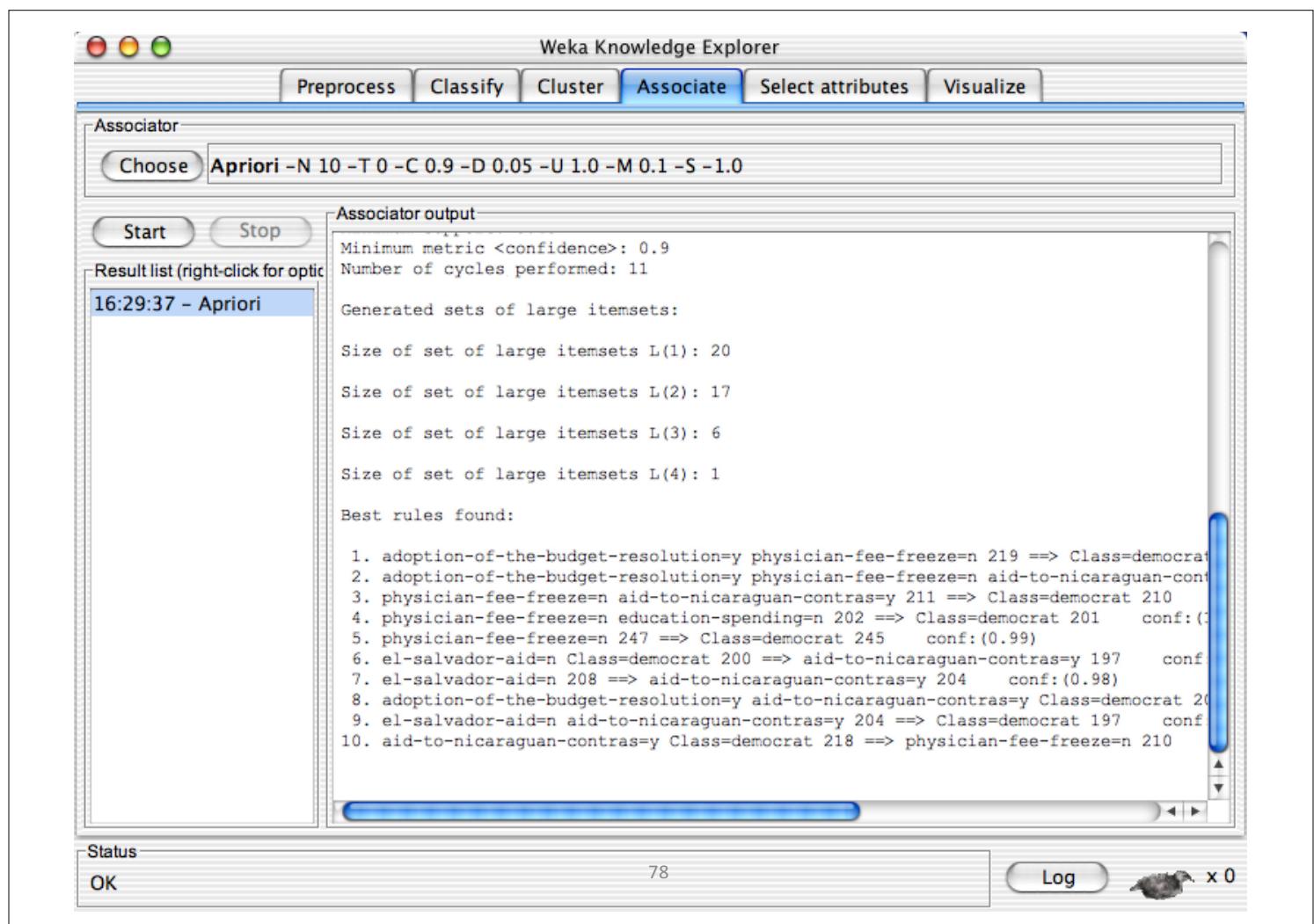
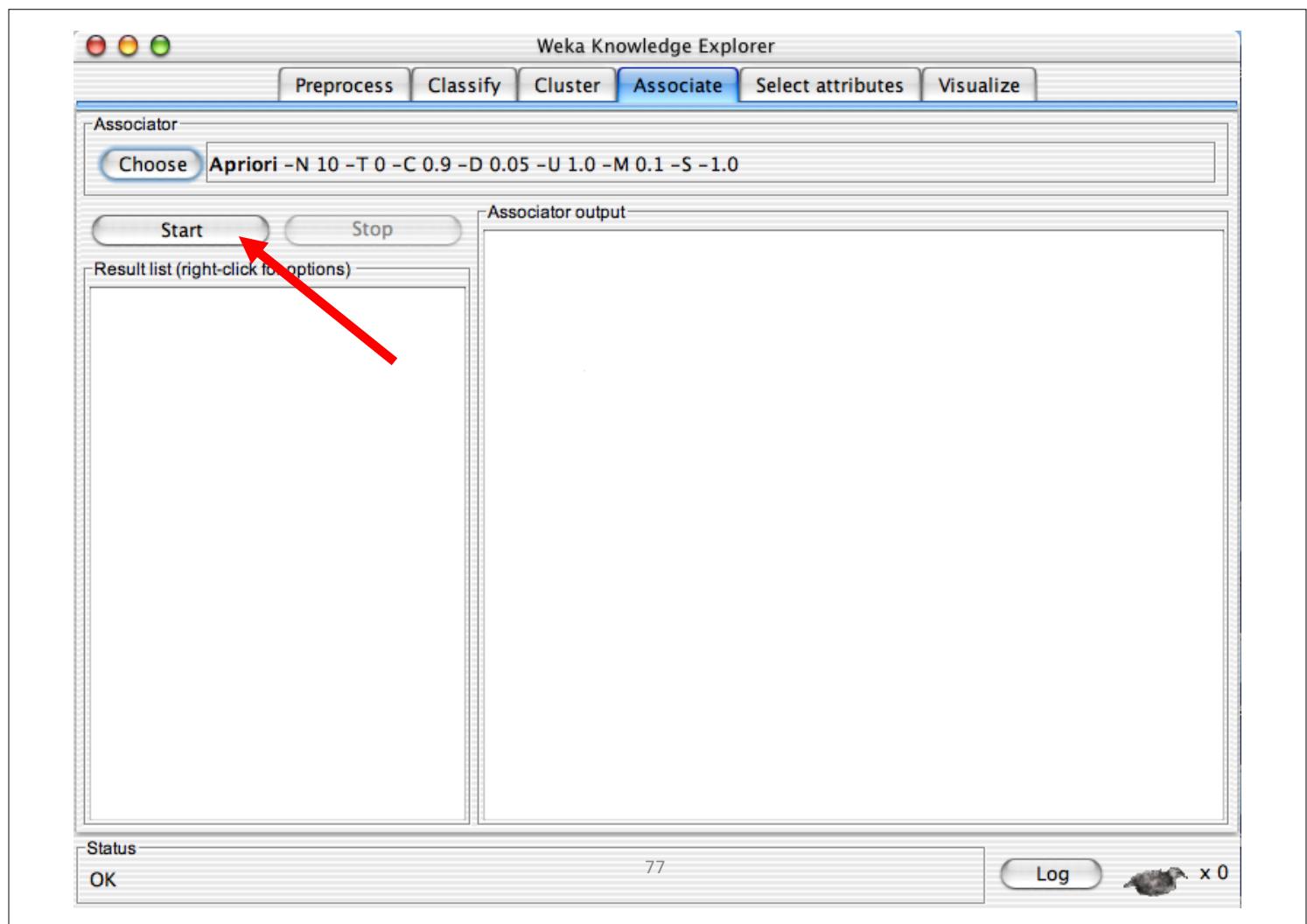
- WEKA contains an implementation of the Apriori algorithm for learning association rules
 - Works only with discrete data
- Can identify statistical dependencies between groups of attributes:
 - milk, butter -> bread, eggs (with confidence 0.9 and support 2000)
- Apriori can compute all rules that have a given minimum support and exceed a given confidence

71





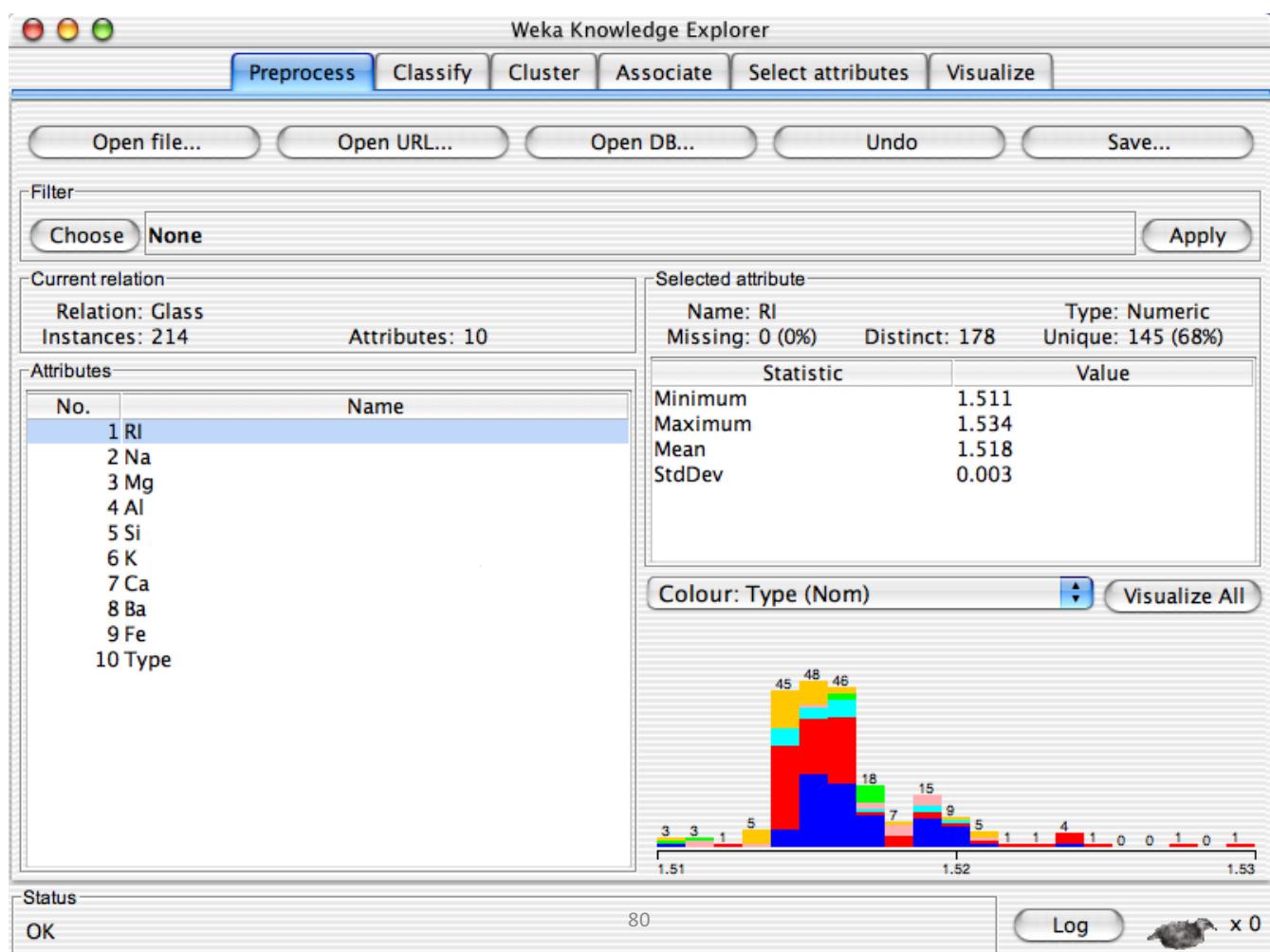


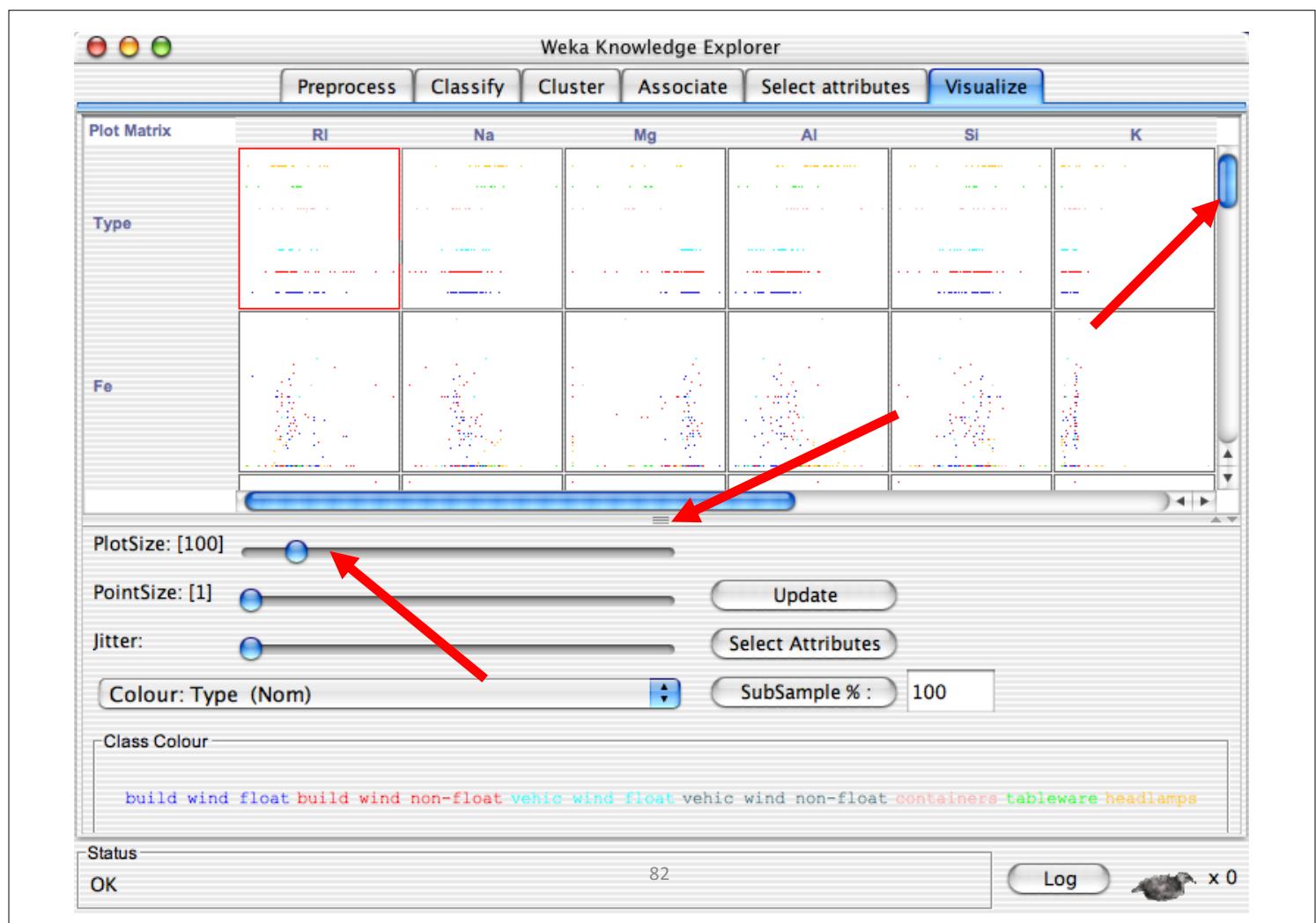
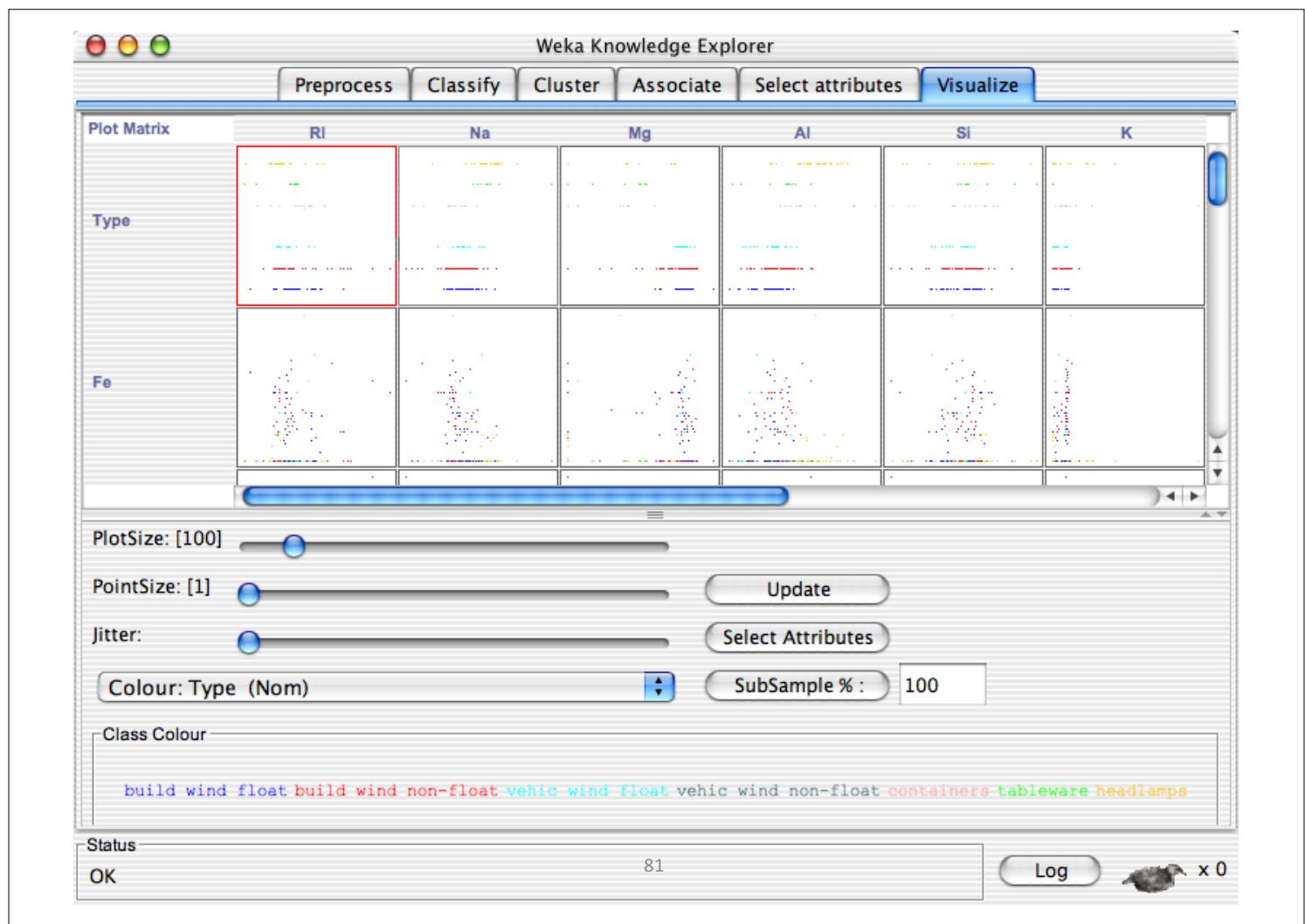


Data visualization

- Visualization very useful in practice:
 - e.g. helps to determine difficulty of the learning problem
- WEKA can visualize single attributes and pairs of attributes
 - To do: rotating 3-d visualizations (Xgobi-style)
- Color-coded class values
- “Jitter” option to deal with nominal attributes (and to detect “hidden” data points)
- “Zoom-in” function

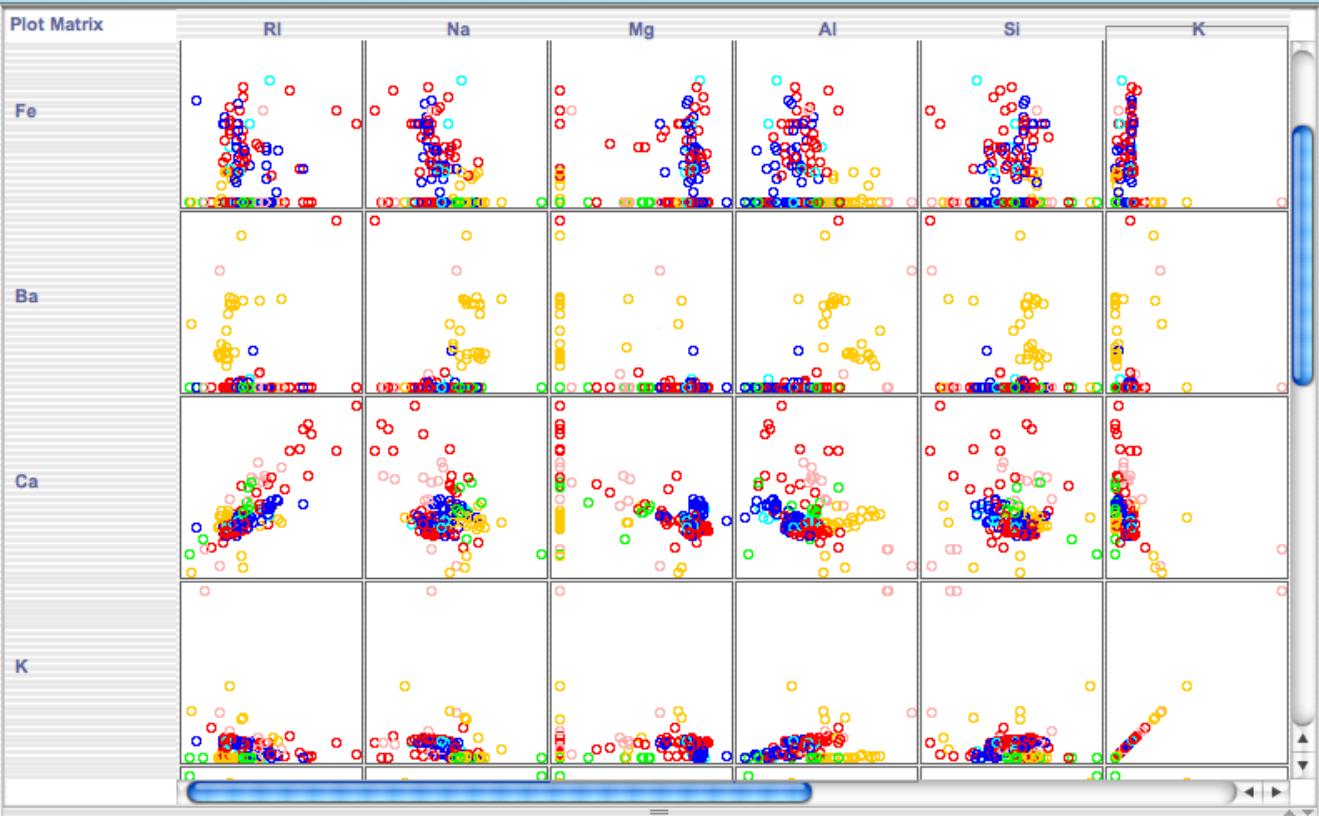
79





Weka Knowledge Explorer

Preprocess Classify Cluster Associate Select attributes Visualize



Status

OK

83

Log

