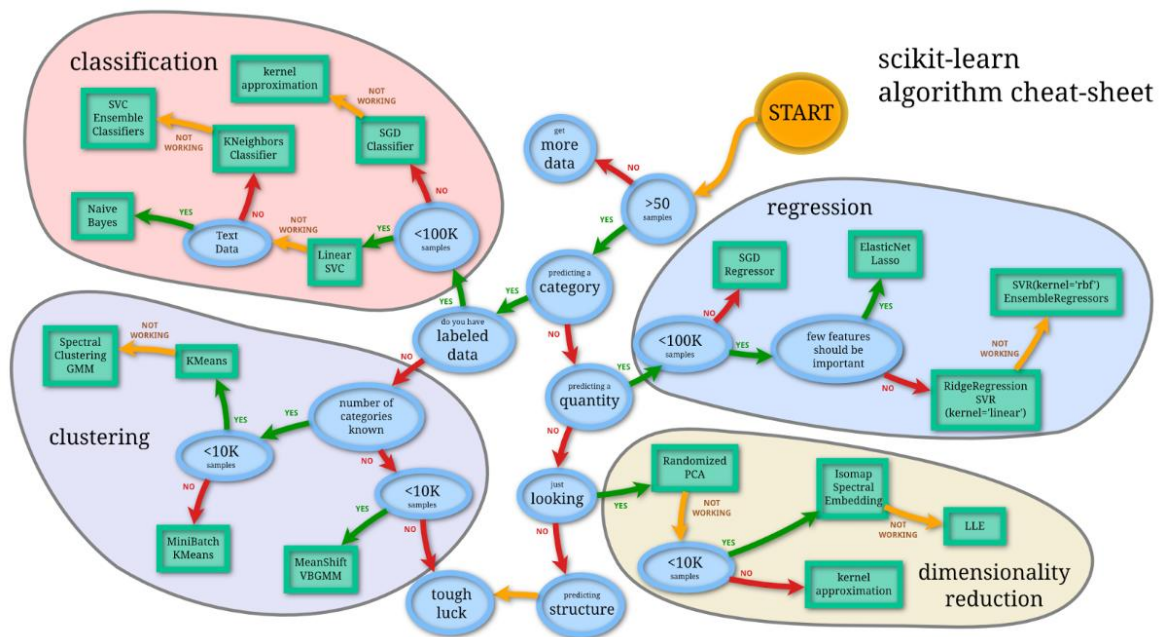# Overview: Data Mining Methods



# 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

---

# 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, …

# 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
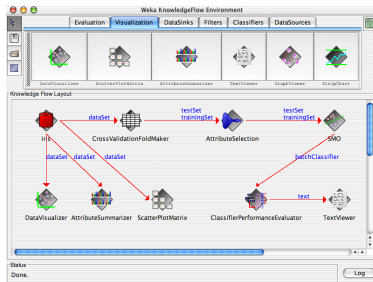67,male,asympt,286,yes,present
67,male,asympt,229,yes,present
38,female,non_anginal,?,no,not_present
…

*Flat file in ARFF format*

---

# WEKA with "flat" files

@relation heart-disease-simplified

*numeric attribute*
*nominal attribute*

@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
…

Weka Knowledge Flow Environment

Evaluation | Visualization | DataSinks | Filters | Classifiers | DataSources

Knowledge Flow Layout

iris   CrossValidationFoldMaker   AttributeSelection   SMO

DataVisualizer   AttributeSummarizer   ScatterPlotMatrix   ClassifierPerformanceEvaluator   TextViewer

Status
Done.   Log

Weka Experiment Environment

Setup | Run | Analyse

Experiment Configuration Mode:   ○ Simple   ○ Advanced

Open...   Save...   New

Results Destination
JDBC database   URL: jdbc:idb=experiments.prp   Browse...

Experiment Type
Cross-validation
Number of folds: 10
● Classification   ○ Regression

Iteration Control
Number of repetitions: 10
● Data sets first
○ Algorithms first

Datasets
Add new...   Delete selected
☐ Use relative paths
/Users/eibe/Documents/datasets/UCI/iris.arff
/Users/eibe/Documents/datasets/UCI/vote.arff
/Users/eibe/Documents/datasets/UCI/glass.arff

Algorithms
Add new...   Delete selected
J48 -C 0.25 -M 2
NeuralNetwork -L 0.3 -M 0.2 -N 500 -V 0 -S 0 -E 20 -H a
NaiveBayes

Notes

Weka GUI Chooser

Visualization   Tools   Help

...vironment for Knowledge Analysis
...10
...013
...ity of Waikato
...ew Zealand

WEKA
The University
of Waikato

Applications

Explorer

Experimenter
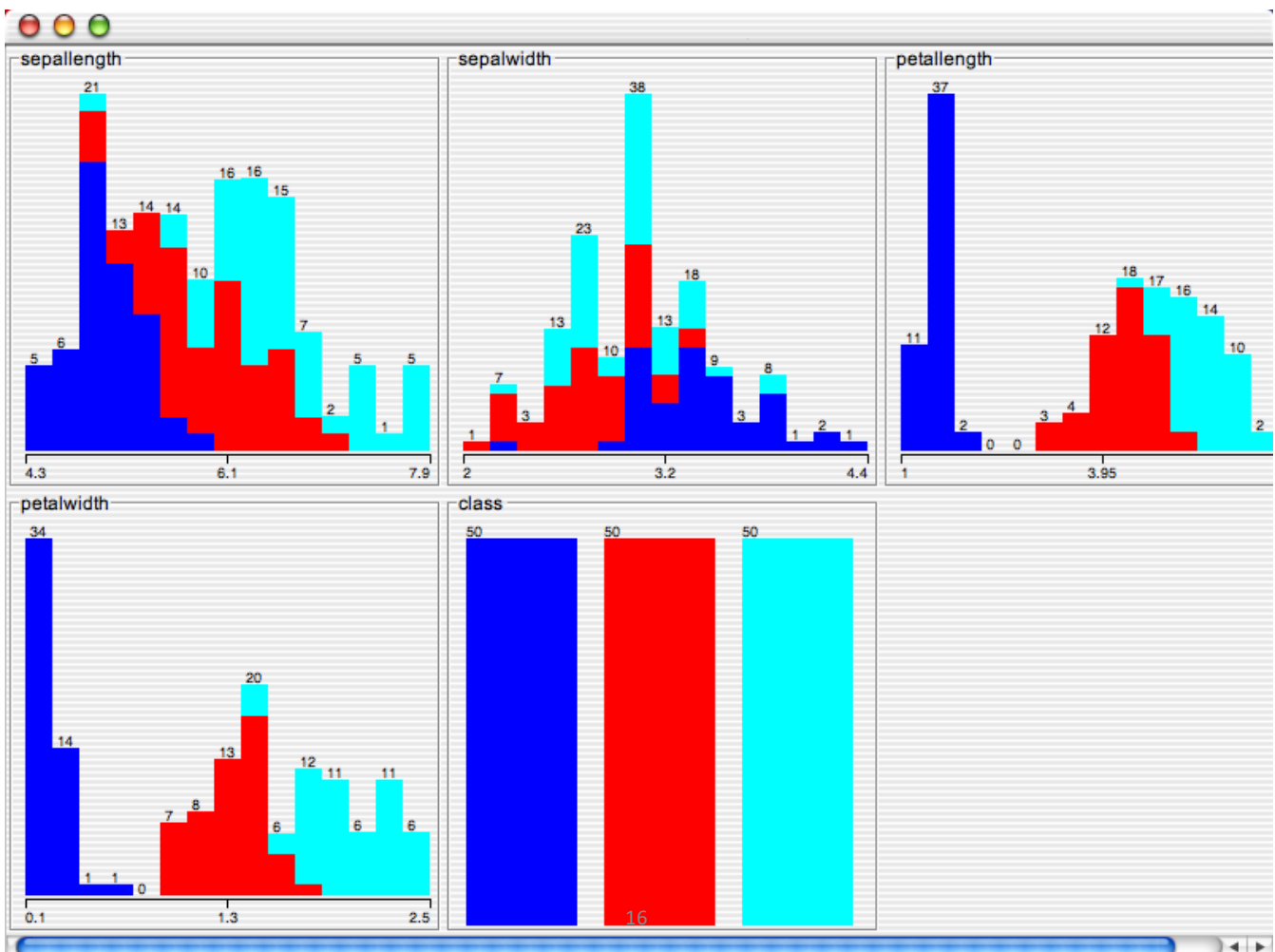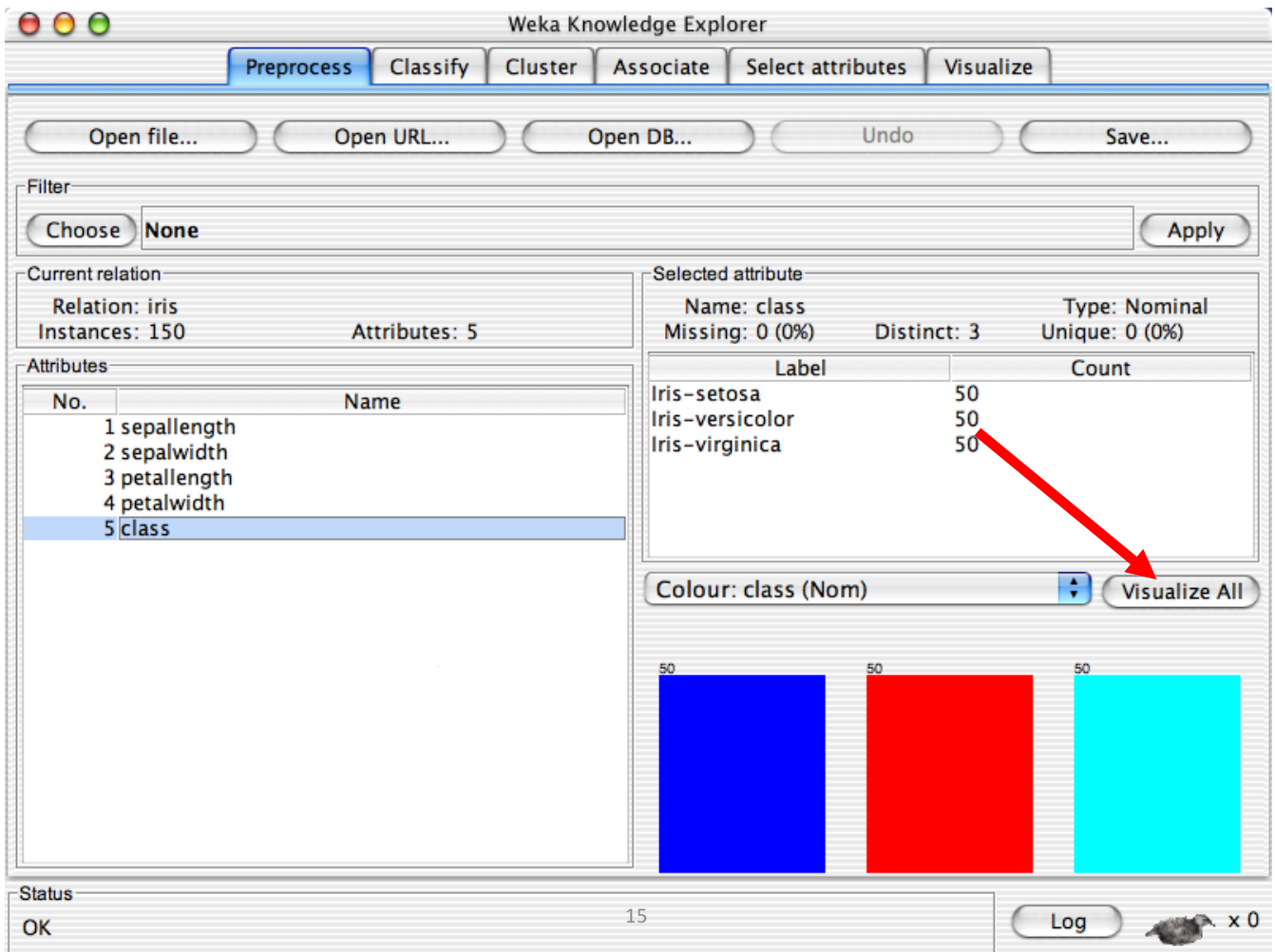
KnowledgeFlow

Simple CLI

Welcome to the WEKA SimpleCLI

Enter commands in the textfield at the bottom of
the window. Use the up and down arrows to move
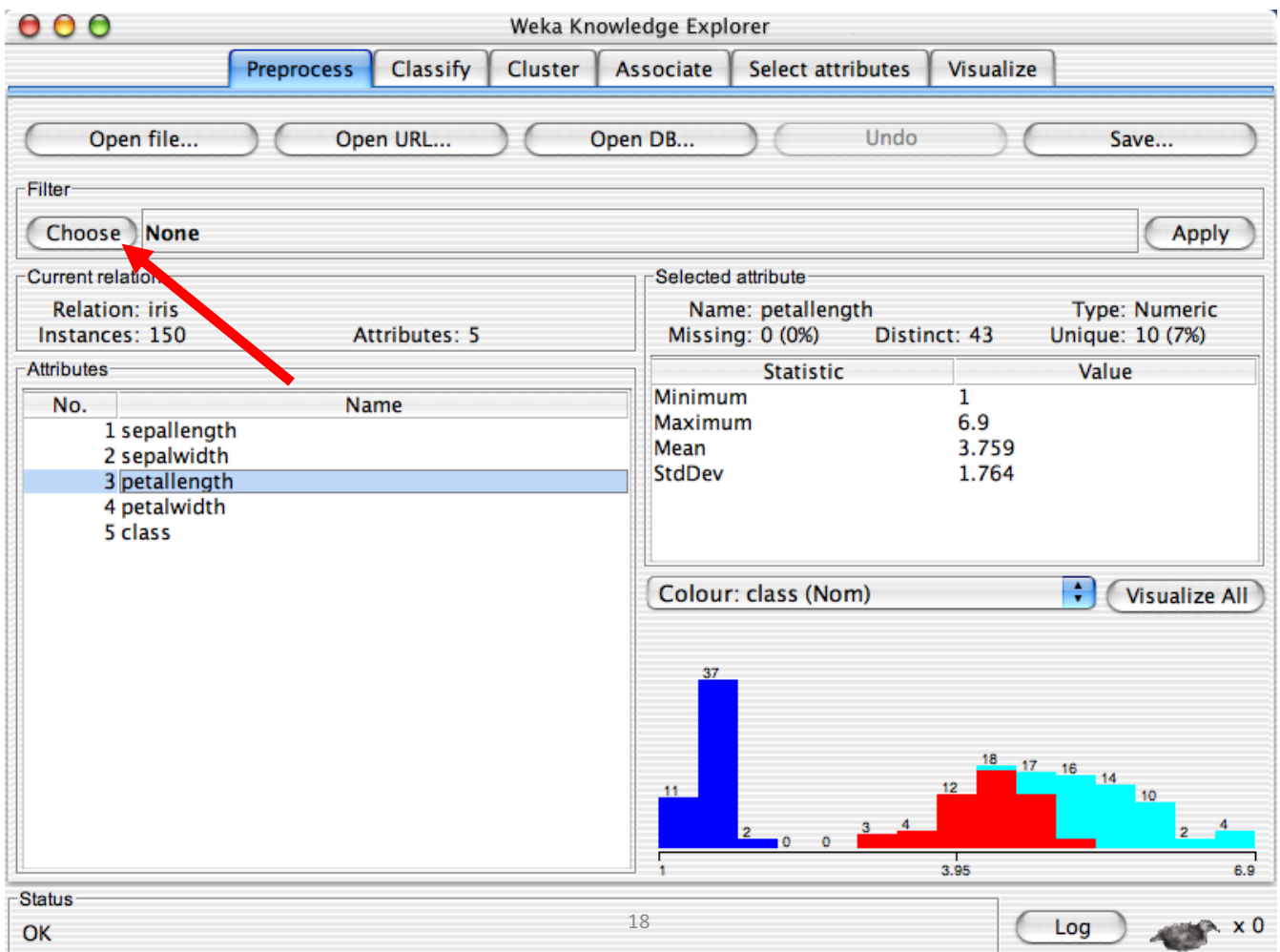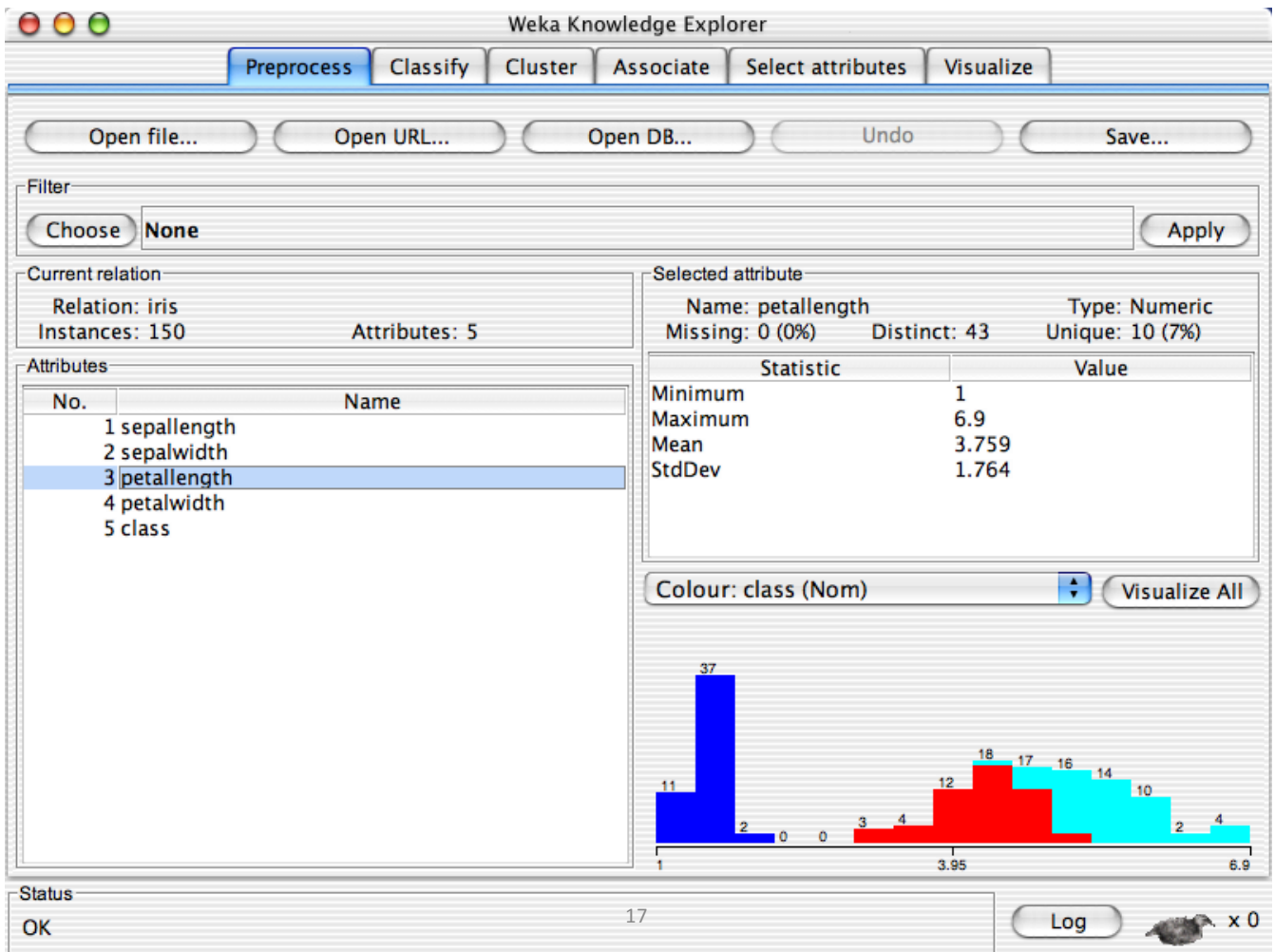through previous commands.

> help

Command must be one of:
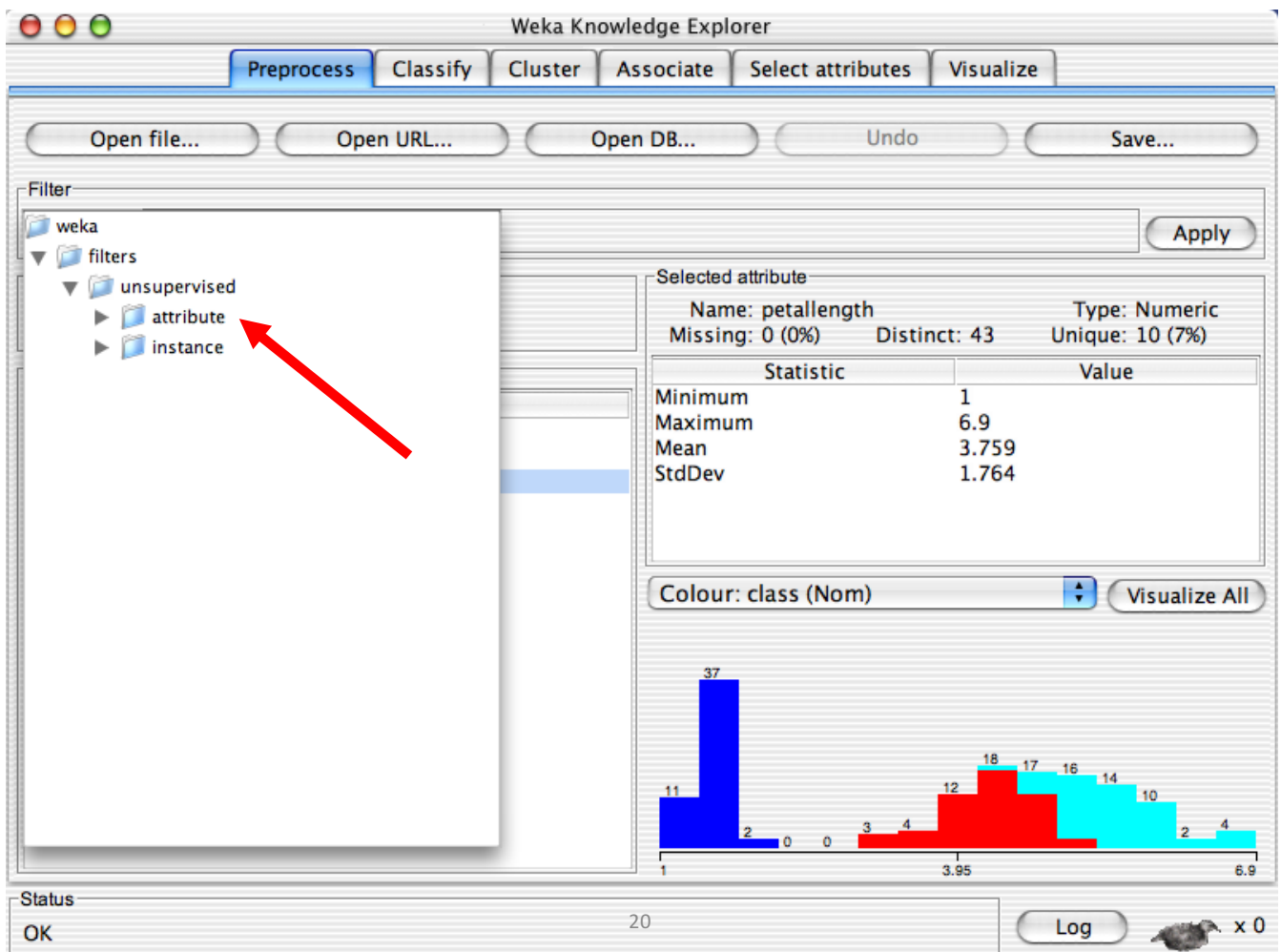    java <classname> <args>
    break
    kill
    cls
    exit
    help <command>

9



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: None
Instances: None   Attributes: None

Selected attribute
Name: None   Type: None
Missing: None   Distinct: None   Unique: None

Attributes

Visualize All

Status
Welcome to the Weka Knowledge Explorer

10

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: None
Instances: None | Attributes: None

Attributes

Selected attribute
Name: None | Type: None
Missing: None | Distinct: None | Unique: None

Visualize All

Status
Welcome to the Weka Knowledge Explorer          11          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: 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

4.3 | 6.1 | 7.9

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

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 |

Attributes

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

Colour: class (Nom) | Visualize All

50    50    50

Status
OK          15          Log   x 0

---

sepallength

21
16  16
15
14  14
13
10
7
6
5                    5   5
2
1
4.3          6.1          7.9

sepalwidth

38
23
18
13        13
10
9
7         8
3                    3
1                1   2   1
2          3.2          4.4

petallength

37
18  17
16
14
12
11                10
4
3
2        0   0          2
1          3.95

petalwidth

34
20
14
13
12
11  11
7   8
6        6   6
1   1
0
0.1          1.3          2.5

class

50    50    50

16

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

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

**Filter**

Apply

- 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

37
11
2   0   0   3   4   12   18   17   16   14   10   2   4
1                    3.95                    6.9

**Status**

OK                                    21                    Log    x 0

---

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

37
11
2   0   0   3   4   12   18   17   16   14   10   2   4
1                    3.95                    6.9

**Status**

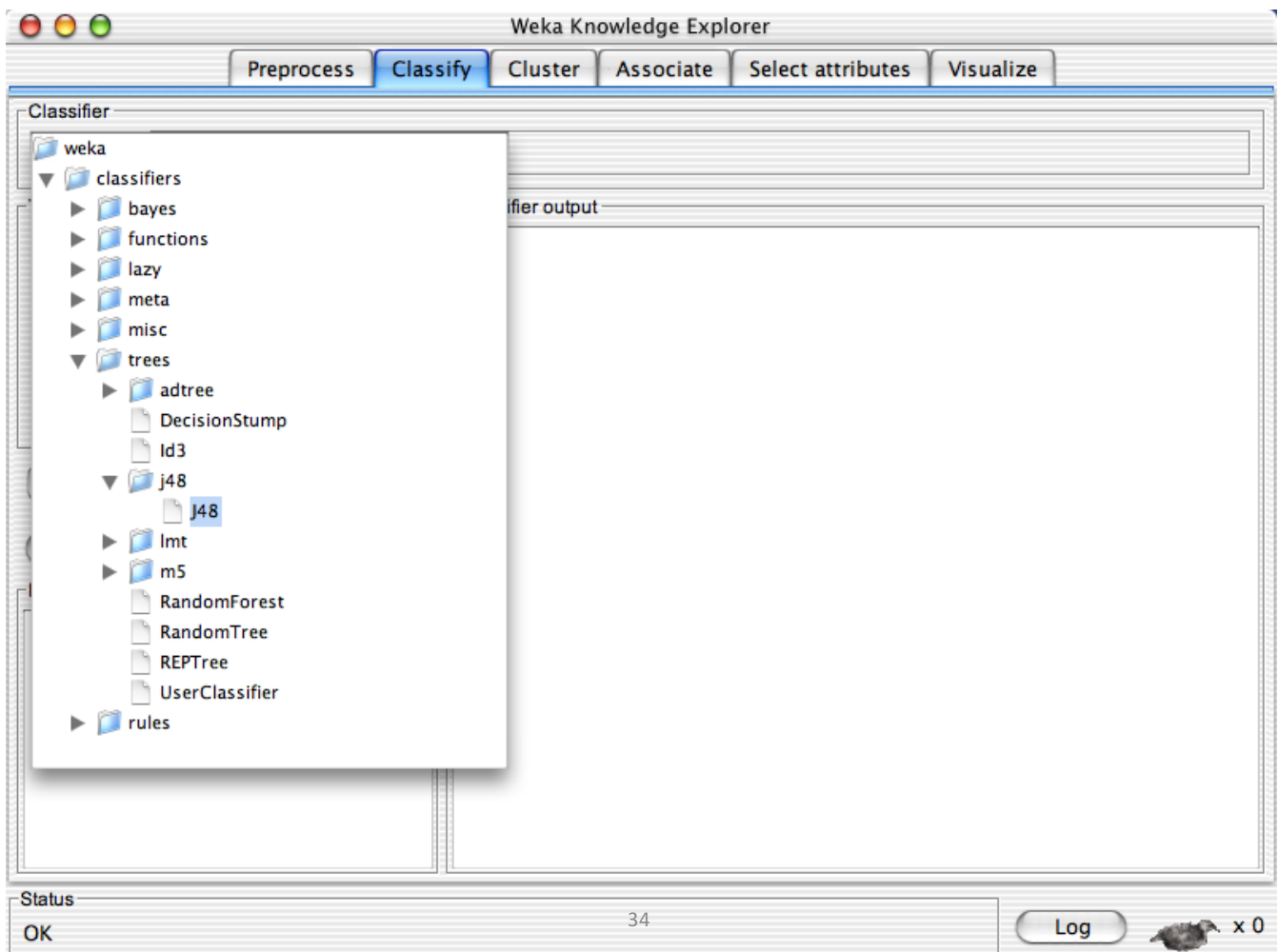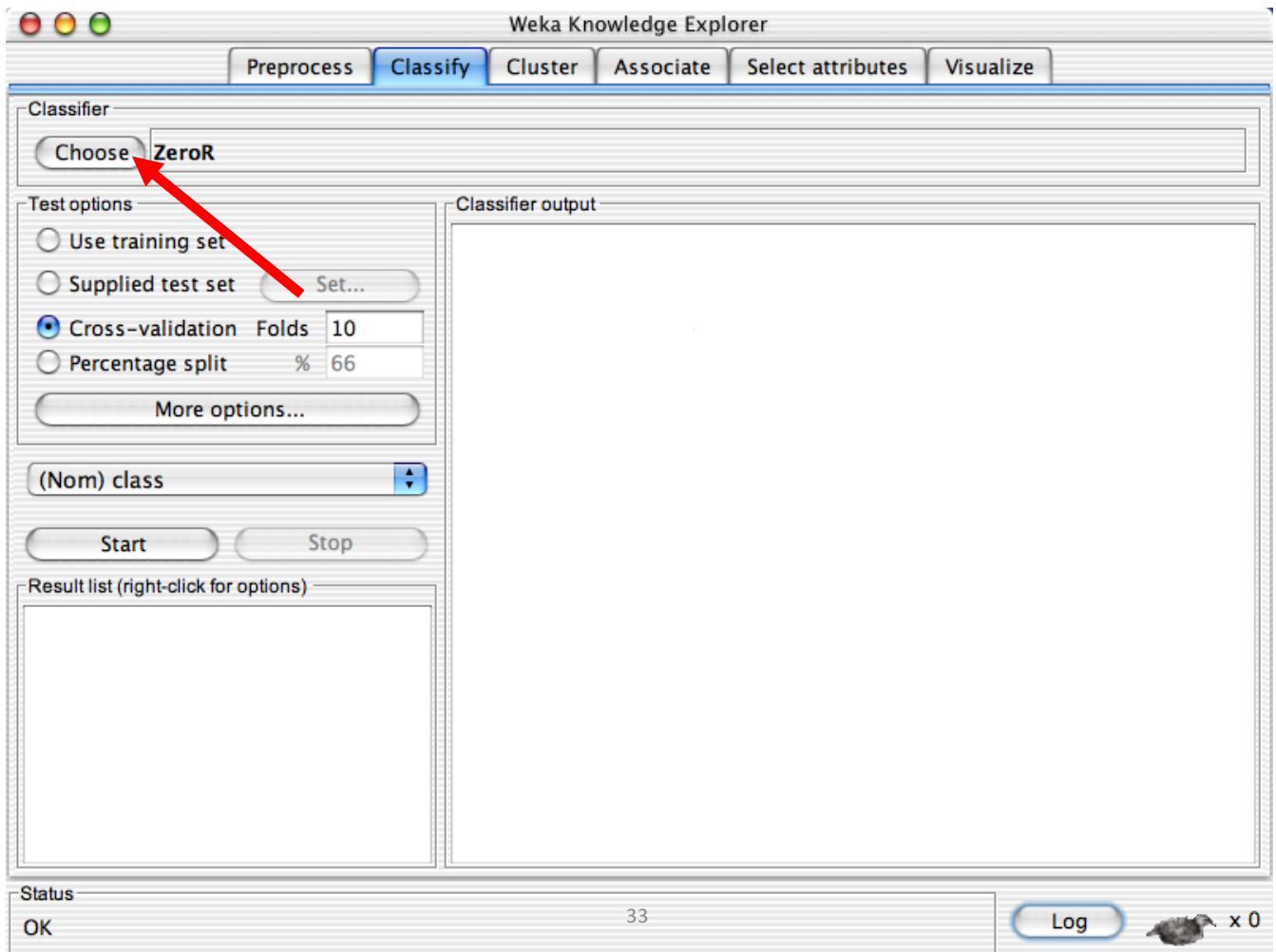OK                                    22                    Log    x 0
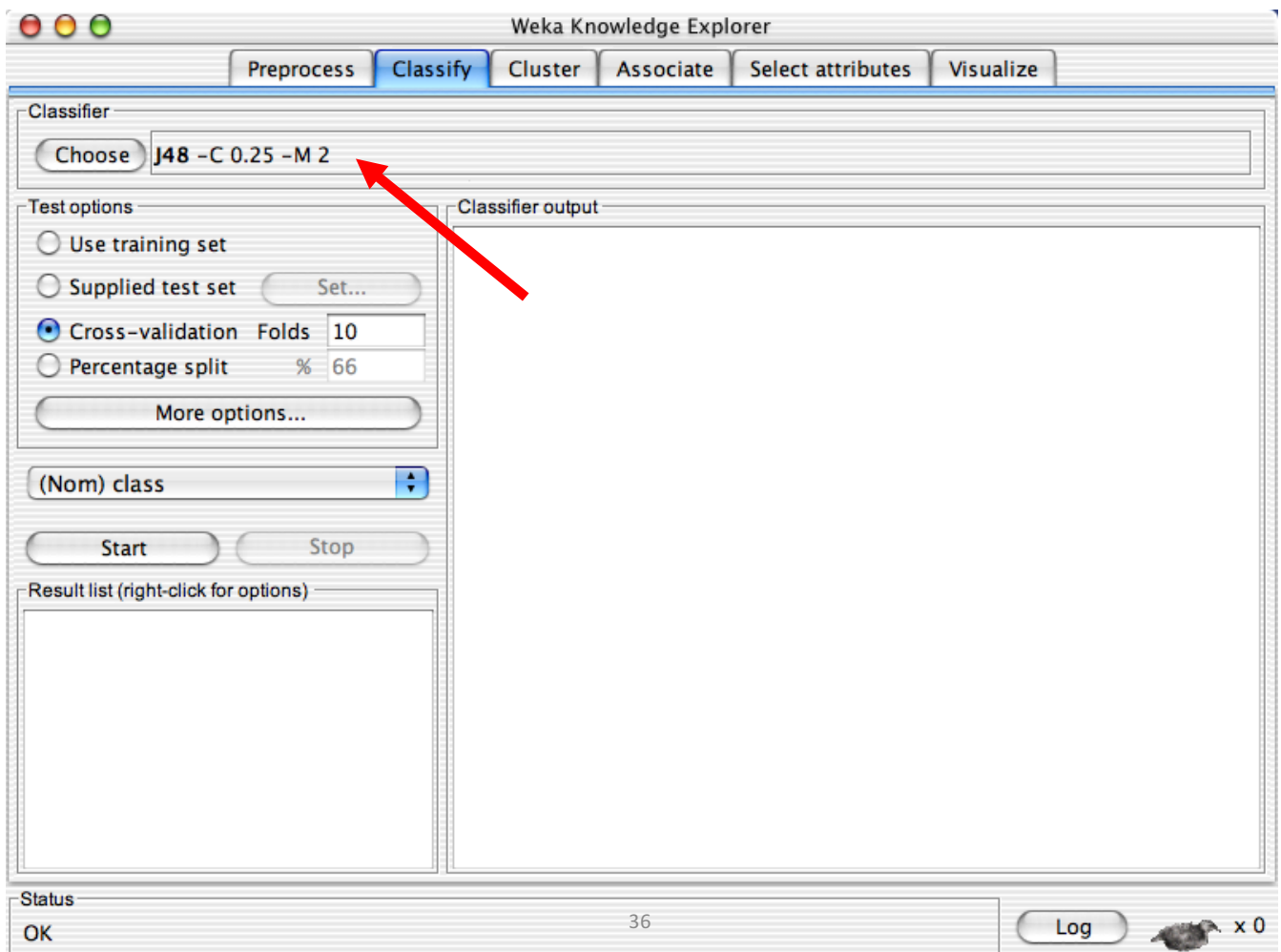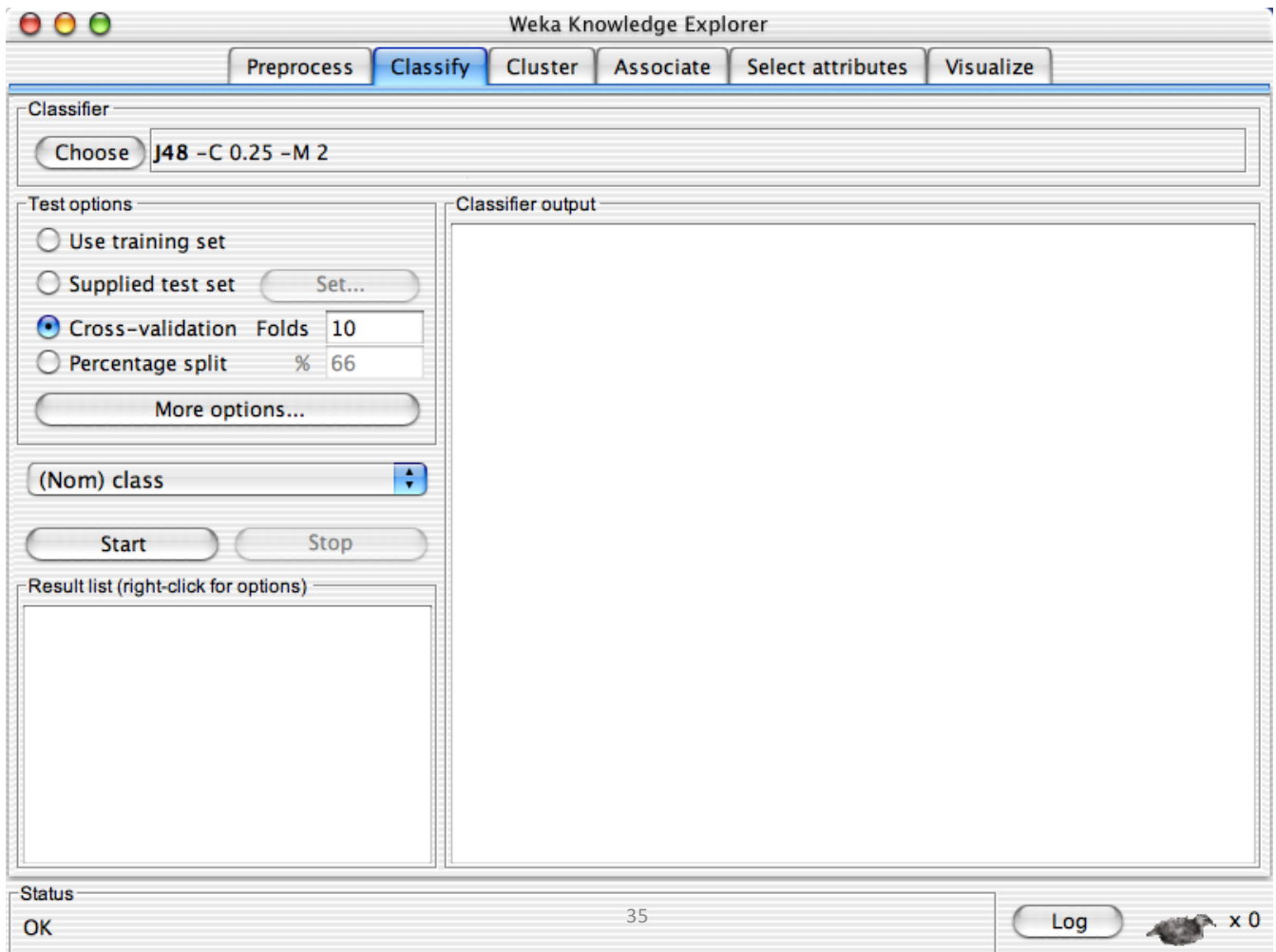
# Building "Classifiers"

- 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:

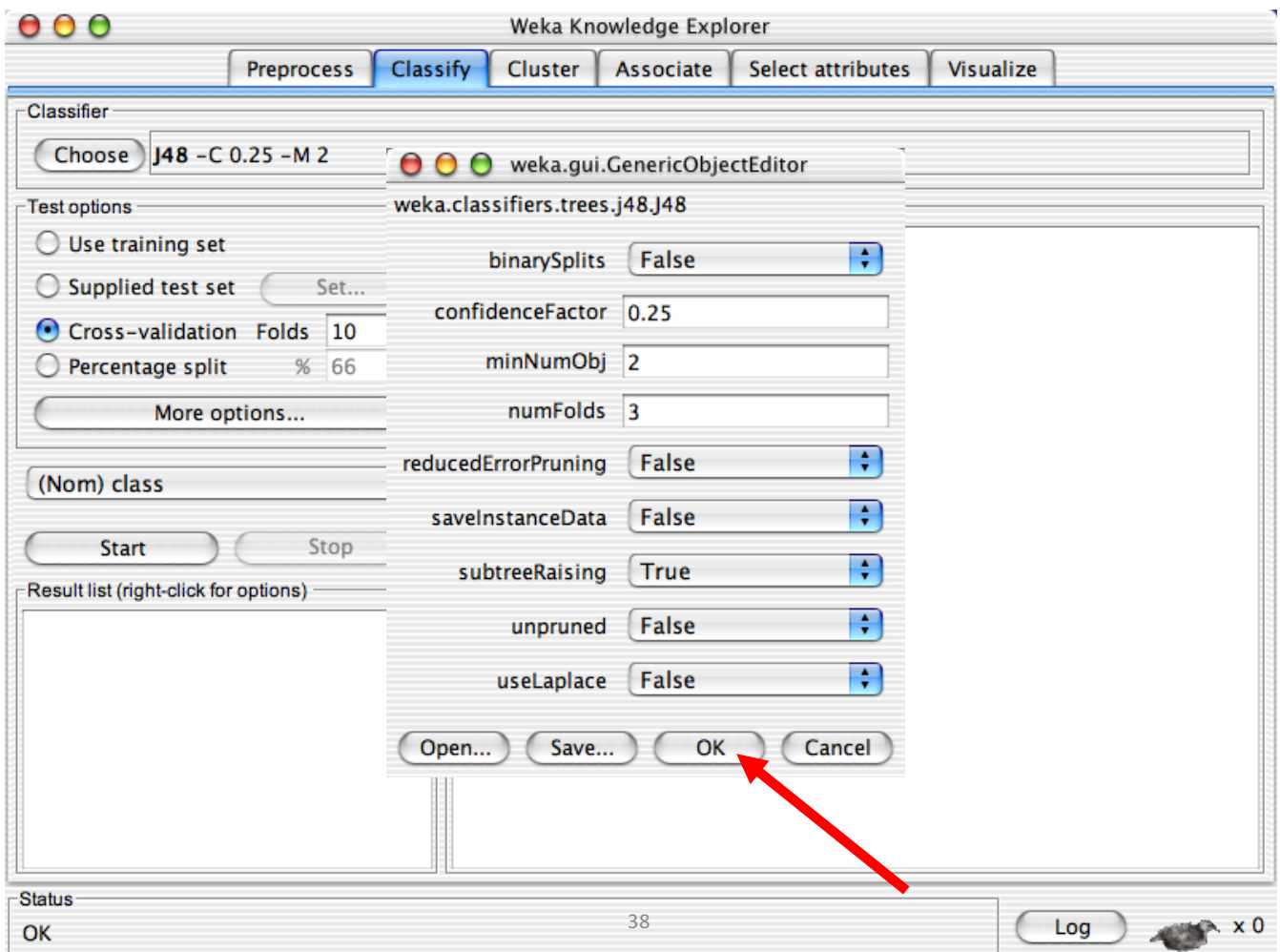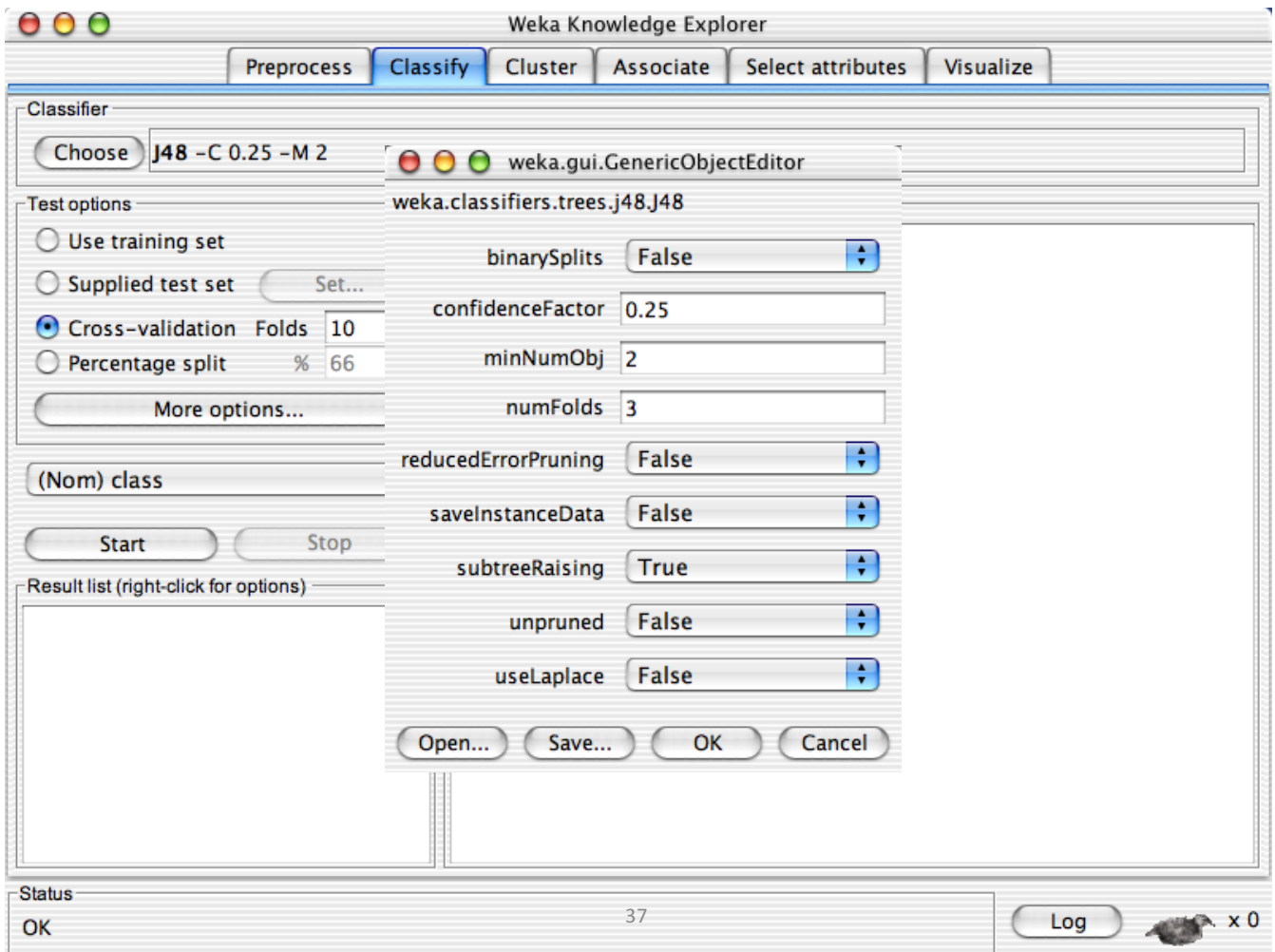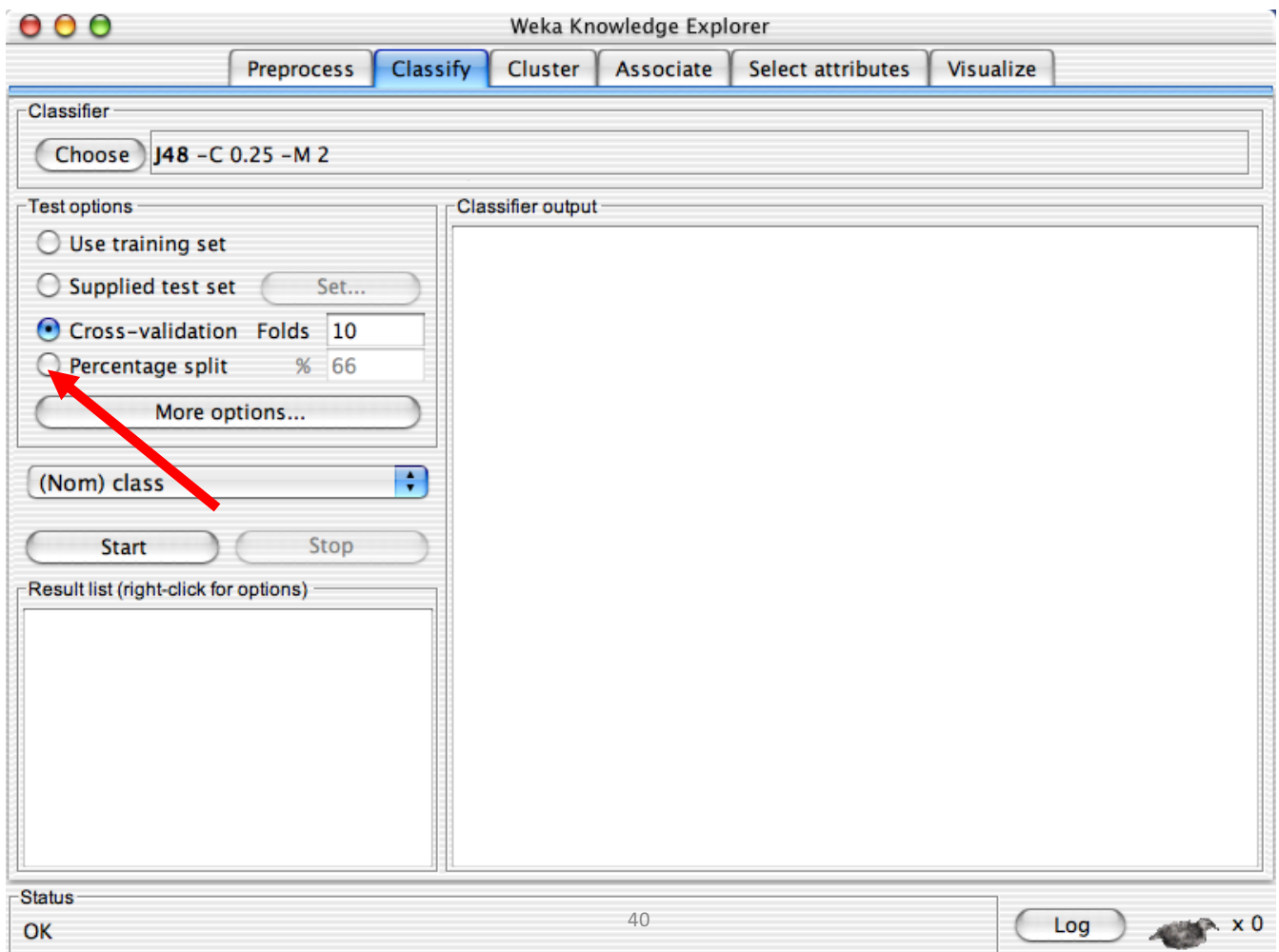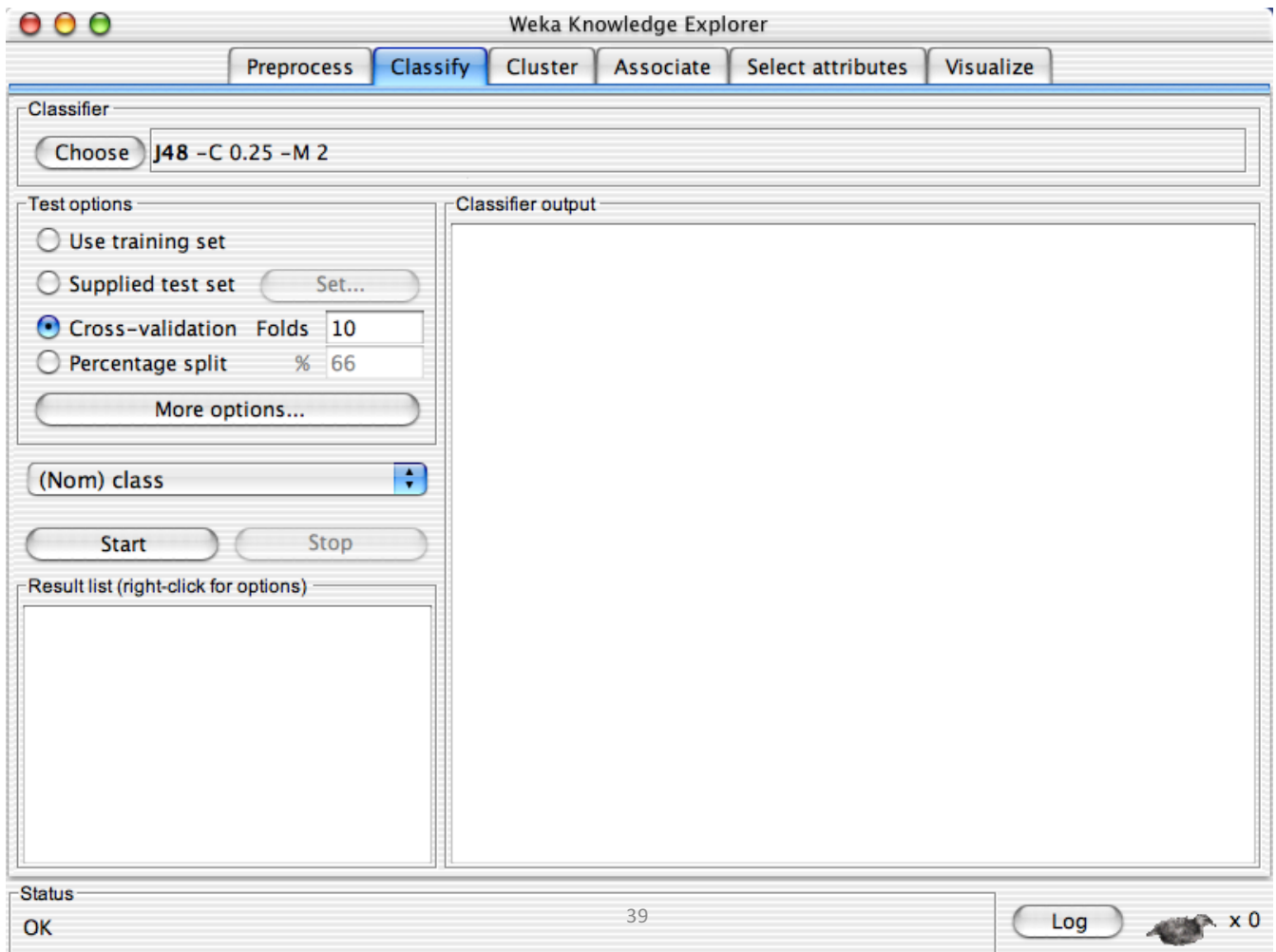  - Bagging, boosting, stacking, error-correcting output codes, locally weighted learning, …

---

=== Run information ===

Scheme:      weka.classifiers.trees.j48.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

**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

**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 ===
```

**Result list (right-click for options)**

11:49:05 – trees.j48.J48

| | | Recall | F-Measure | Class |
|---|---|---|---|---|
| View in main window | | 1 | 1 | Iris-setosa |
| View in separate window | | 1 | 0.95 | Iris-versicolor |
| Save result buffer | | 0.882 | 0.938 | Iris-virginica |
| Load model | | | | |
| Save model | | | | |
| Re-evaluate model on current test set | | | | |
| Visualize classifer 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

**Test options**

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

More opti...

(Nom) class

Start

**Weka Classifier Tree Visualizer: 11:49:05 – trees.j48.J48 (iris)**

**Tree View**



96.0784 %
3.9216 %

**Result list (right-click for**

11:49:05 – trees.j48.J

```
              ass
              is-setosa
              is-versicolor
              is-virginica
```

```
13  0  0 |  a = Iris-setosa
 0 19  0 |  b = Iris-versicolor
 0  2 15 |  c = Iris-virginica
```

**Status**

OK    52

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

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
```
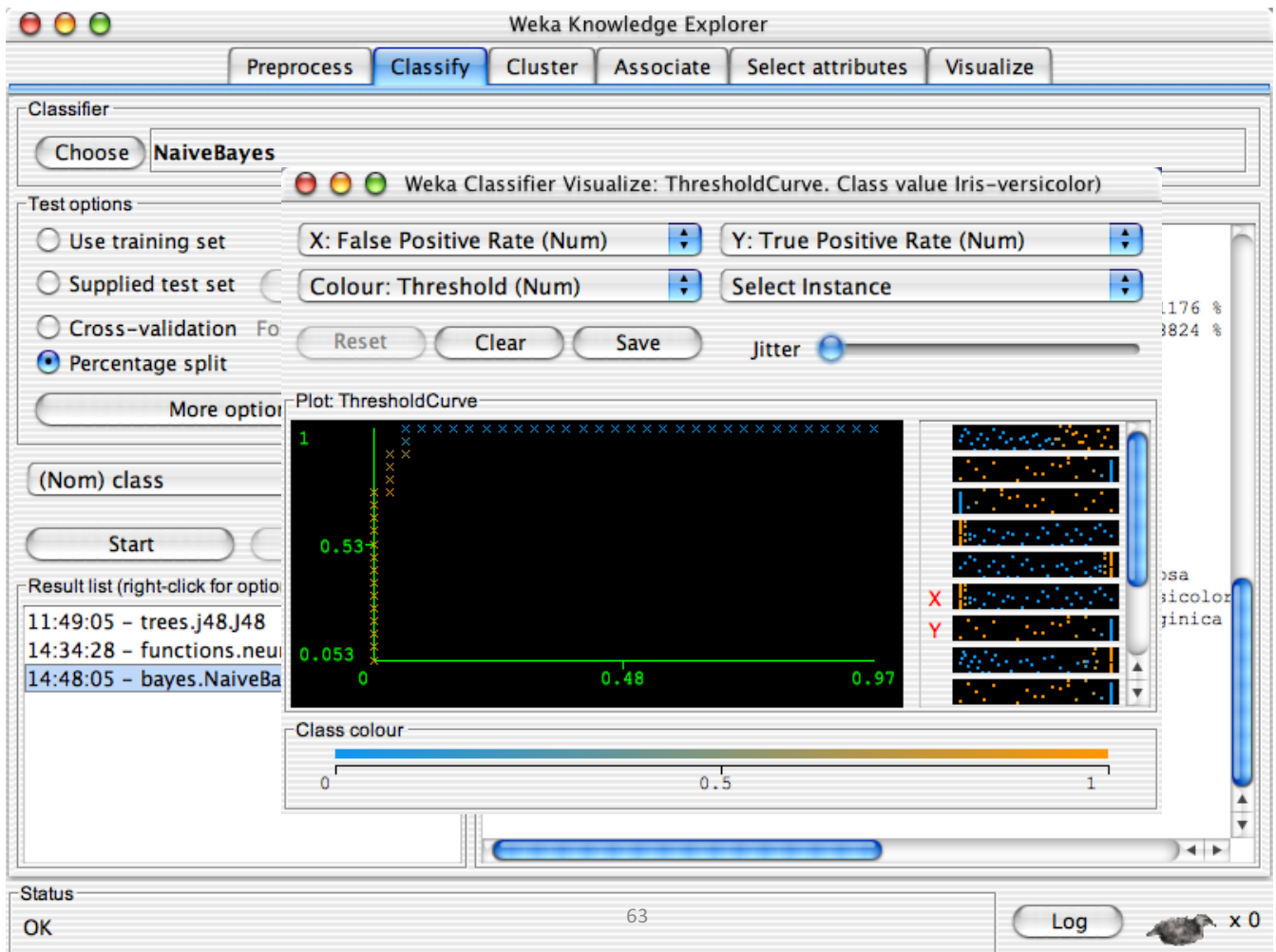
Status
OK

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 ===

orrectly Classified Instances      50            98.0392 %
ncorrectly Classified Instances     1             1.9608 %
appa statistic                     0.9704
ean absolute error                 0.0239
oot mean squared error             0.1101
elative absolute error             5.3594 %
oot relative squared error        23.2952 %
otal Number of Instances           51

== Detailed Accuracy By Class ===

P 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

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
```
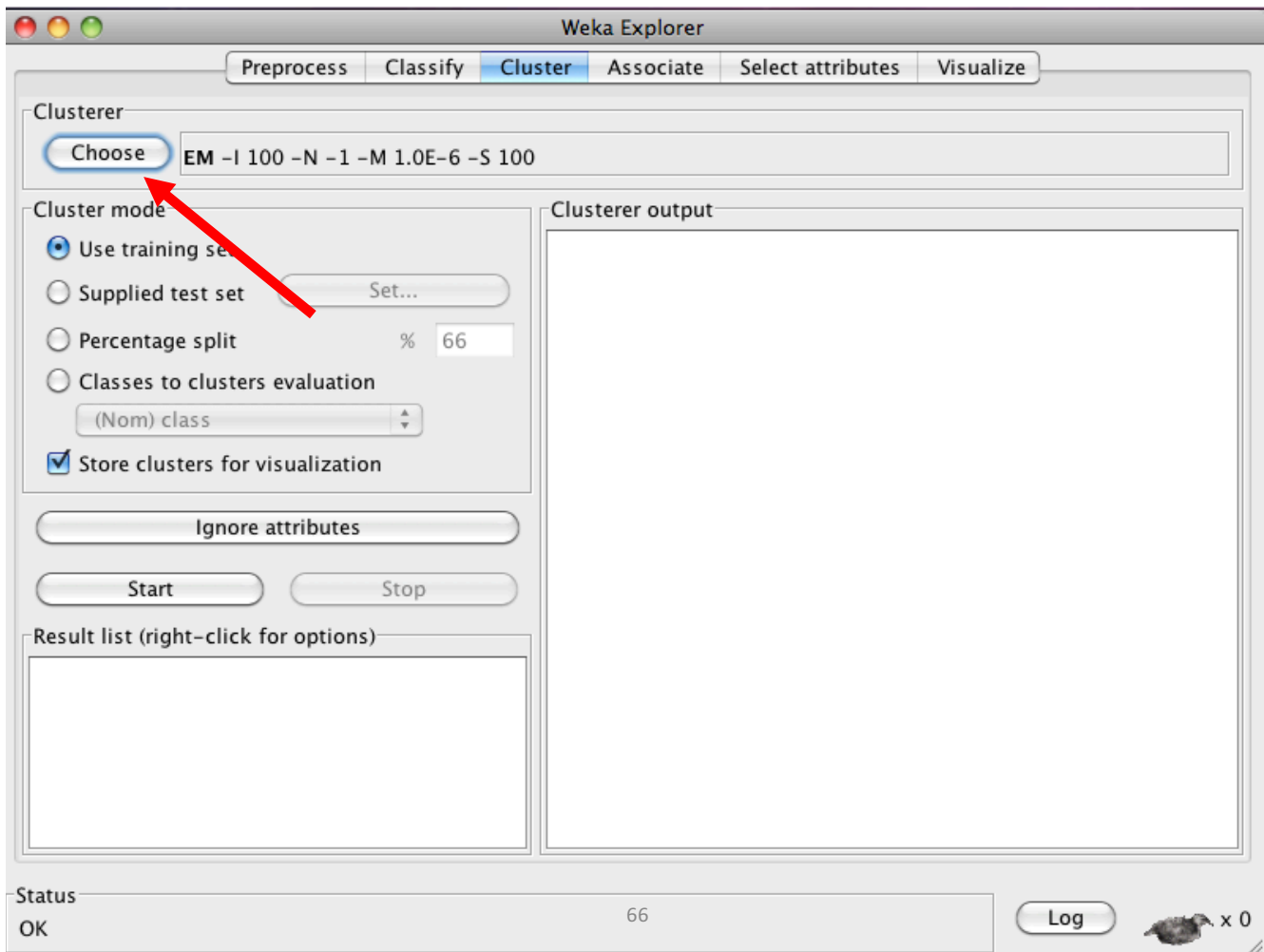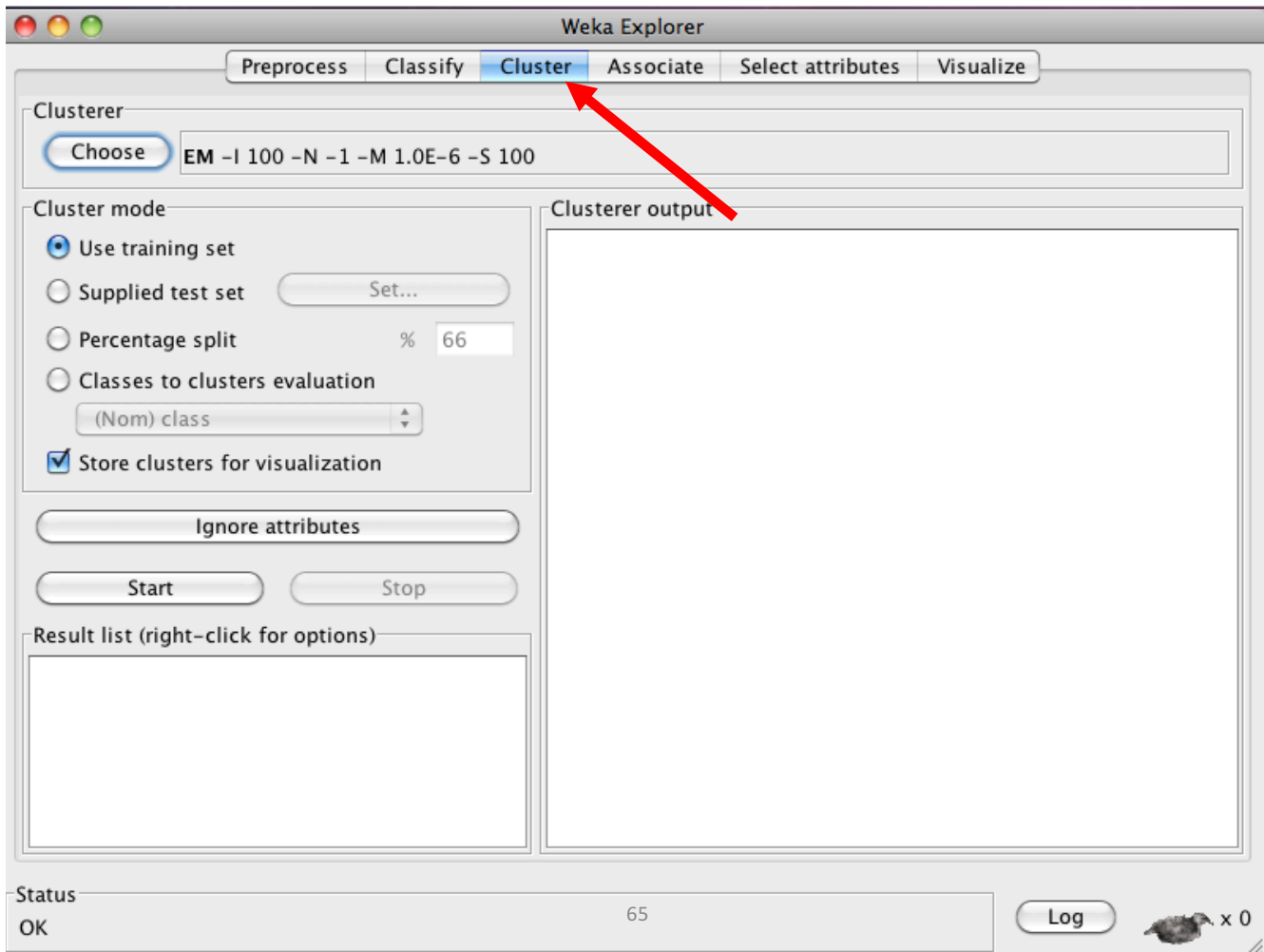
**Status**

Problem evaluating classifier

Log   x 0

# **Clustering**

- WEKA contains many clustering implementations:
  - Works with both discrete and numerical data

- Example of K-means

# 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

---



72

## Weka Knowledge Explorer

Preprocess  Classify  Cluster  **Associate**  Select attributes  Visualize

Associator

Choose | **Apriori** -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0

Start  Stop

Associator output

Result list (right-click for options)

Status

OK

77

Log

x 0

---

## Weka Knowledge Explorer

Preprocess  Classify  Cluster  **Associate**  Select attributes  Visualize

Associator

Choose | **Apriori** -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0

Start  Stop

Associator output

Result list (right-click for optic

16:29:37 - Apriori

```
Minimum metric <confidence>: 0.9
Number of cycles performed: 11

Generated sets of large itemsets:

Size of set of large itemsets L(1): 20

Size of set of large itemsets L(2): 17

Size of set of large itemsets L(3): 6

Size of set of large itemsets L(4): 1

Best rules found:

 1. adoption-of-the-budget-resolution=y physician-fee-freeze=n 219 ==> Class=democrat
 2. adoption-of-the-budget-resolution=y physician-fee-freeze=n aid-to-nicaraguan-cont
 3. physician-fee-freeze=n aid-to-nicaraguan-contras=y 211 ==> Class=democrat 210
 4. physician-fee-freeze=n education-spending=n 202 ==> Class=democrat 201    conf:(
 5. physician-fee-freeze=n 247 ==> Class=democrat 245    conf:(0.99)
 6. el-salvador-aid=n Class=democrat 200 ==> aid-to-nicaraguan-contras=y 197    conf
 7. el-salvador-aid=n 208 ==> aid-to-nicaraguan-contras=y 204    conf:(0.98)
 8. adoption-of-the-budget-resolution=y aid-to-nicaraguan-contras=y Class=democrat 20
 9. el-salvador-aid=n aid-to-nicaraguan-contras=y 204 ==> Class=democrat 197    conf
10. aid-to-nicaraguan-contras=y Class=democrat 218 ==> physician-fee-freeze=n 210
```
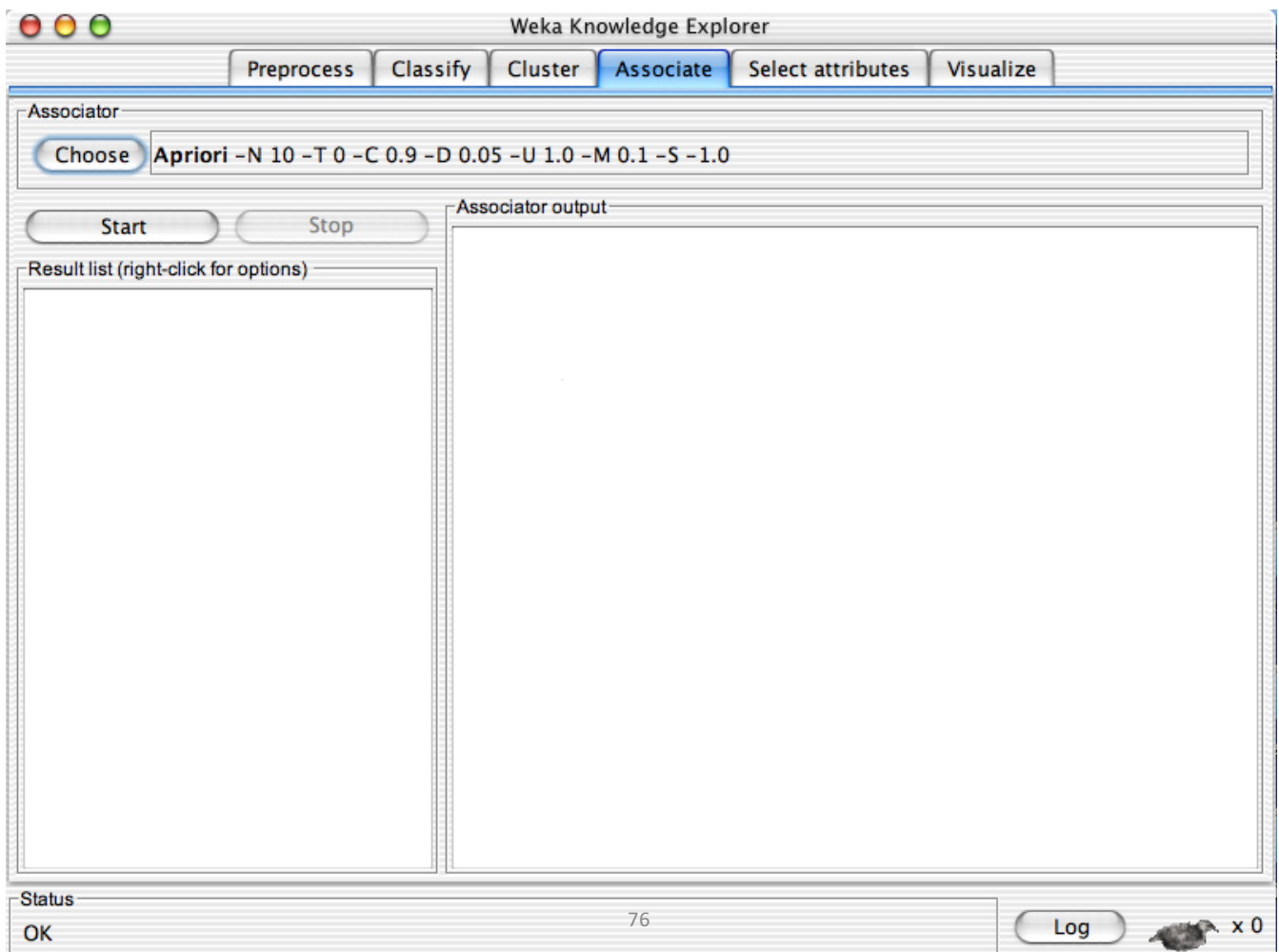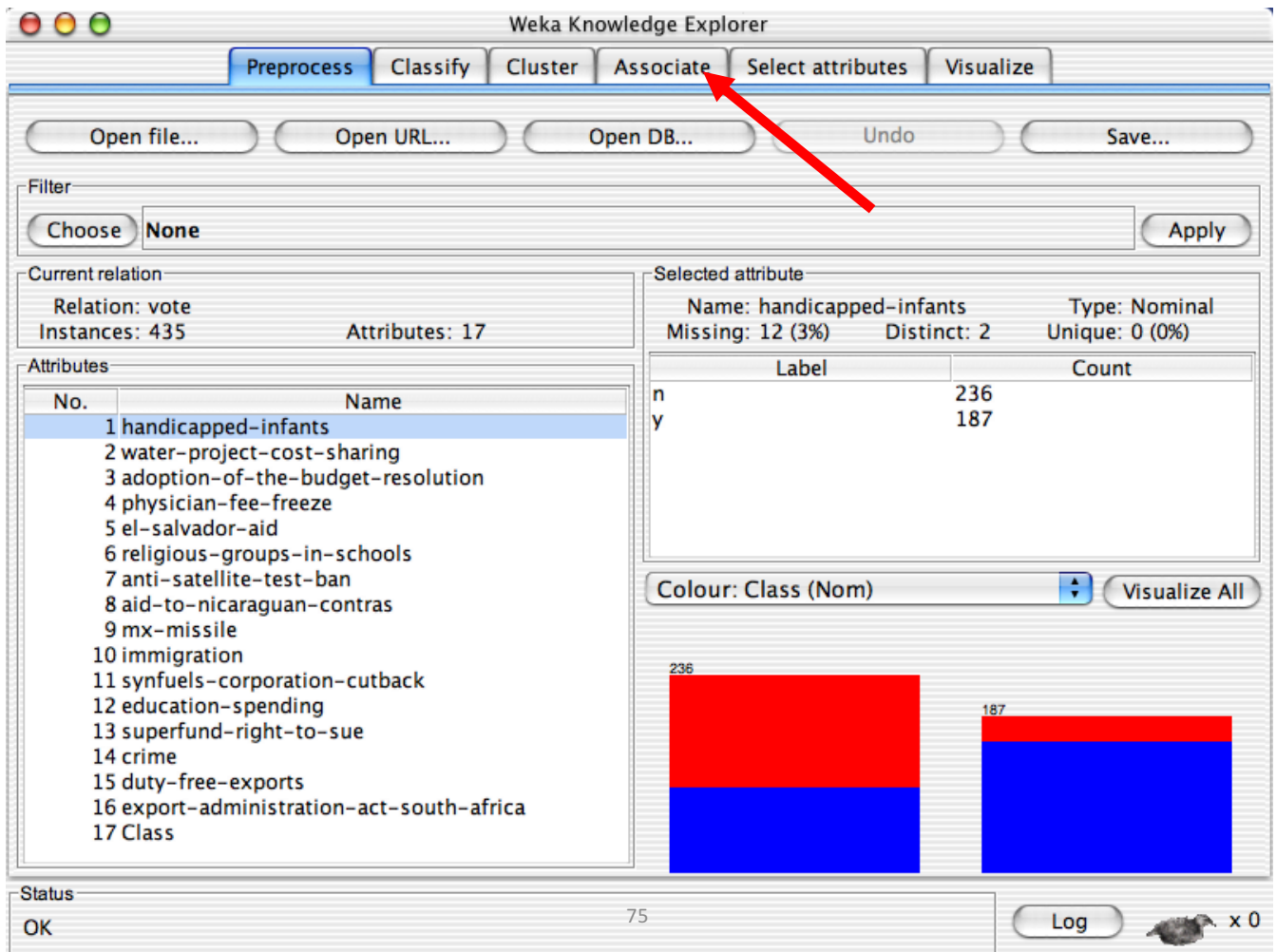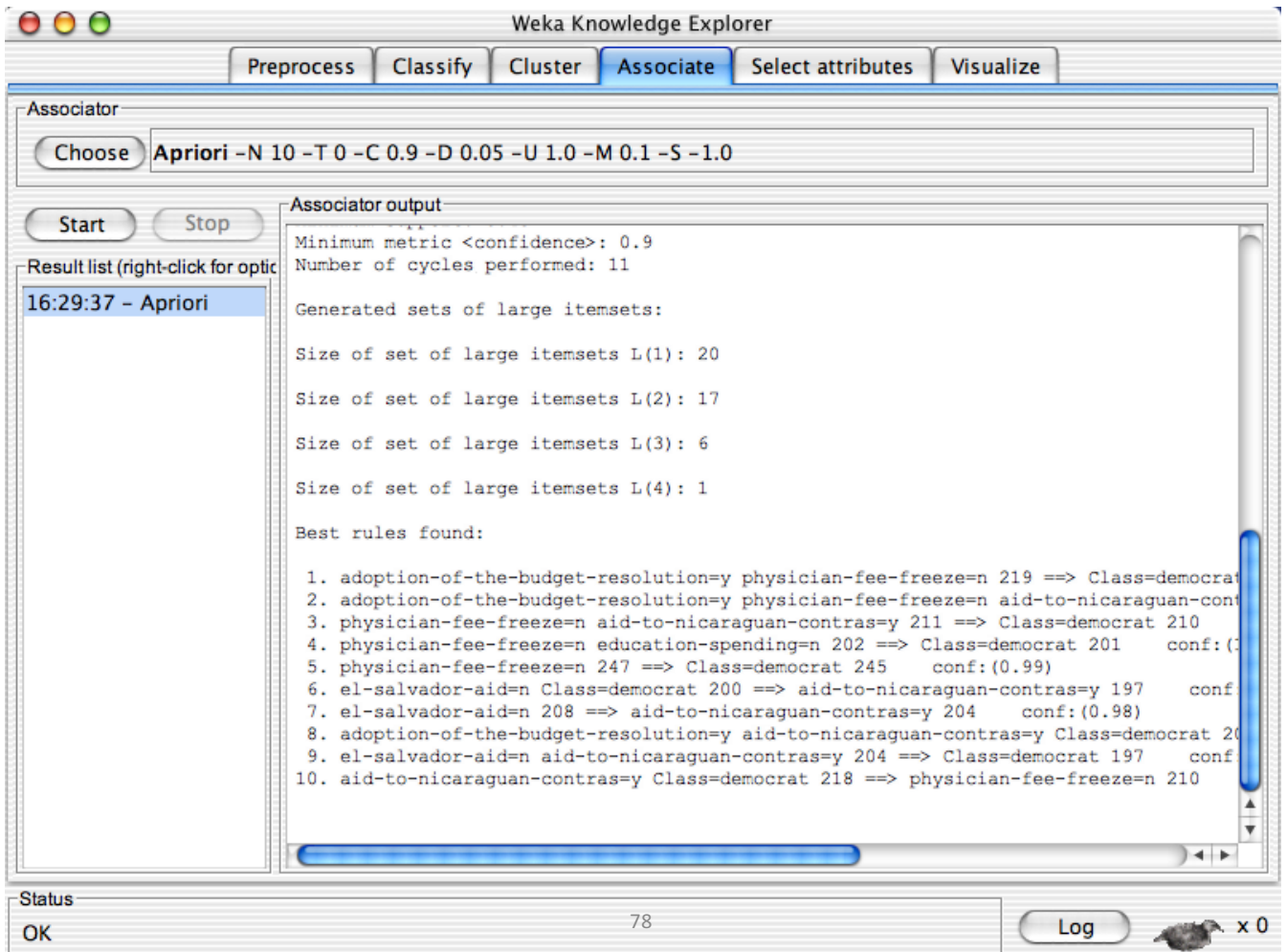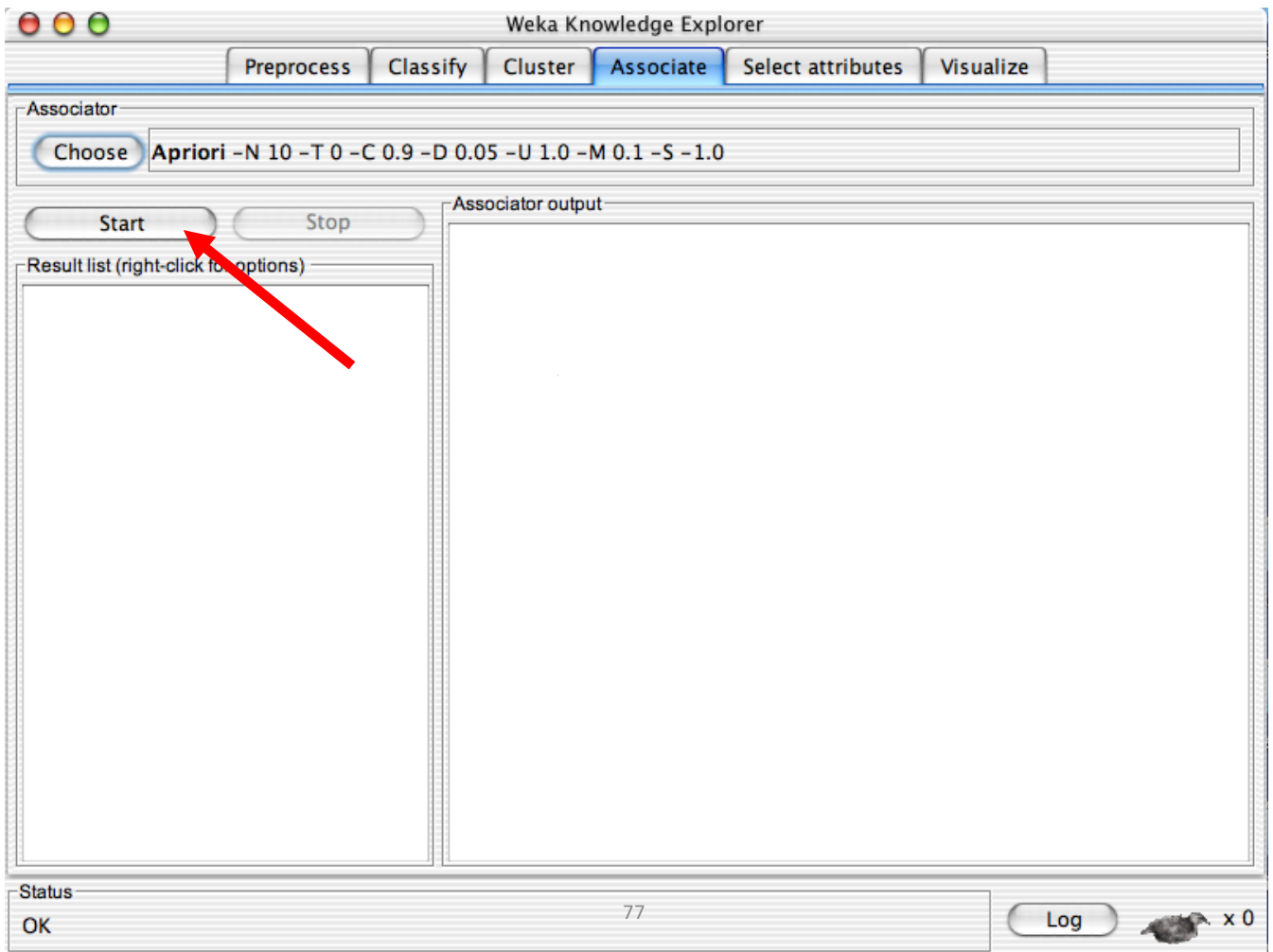
Status

OK

78

Log

x 0

# 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