# Data Mining with WEKA

Original author:  unknown. ? The WEKA team
Additional material: Tim Menzies, 2010

# WEKA

- **Machine learning/data mining software written in Java**
  - Used for research, education, and applications
  - Complements "Data Mining" by Witten & Frank
- **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

# Access

- WEKA is available at

    **http://www.cs.waikato.ac.nz/ml/weka**
- Also has a list of projects based on WEKA
- WEKA contributors:

    Abdelaziz Mahoui, Alexander K. Seewald, Ashraf M. Kibriya, Bern hard Pfahringer , Brent Martin, Peter Flach, Eibe Frank ,Gabi Schmidb erger ,Ian H. Witten , J. Lindgren, Janice Boughton,  Jason Wells, Len Trigg, Lucio de Souza Coelho, Malcolm Ware, Mark Hall ,Remco Bou ckaert , Richard Kirkby, Shane Butler, Shane Legg, Stuart Inglis, Sylva in Roy, Tony Voyle, Xin Xu, Yong Wang, Zhihai Wang
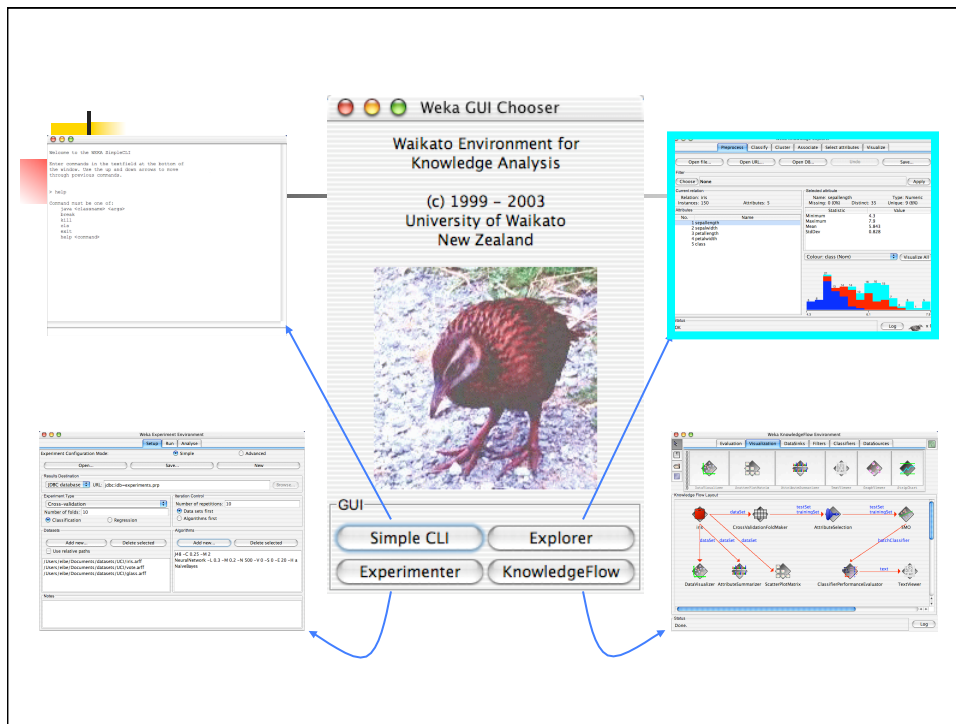
# Data Files

```
@relation heart-disease-simplified

@attribute age numeric
@attribute sex { female, male}
@attribute chest_pain_type { typ_angina, asympt, non_anginal, atyp_angi
    na}
@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
```
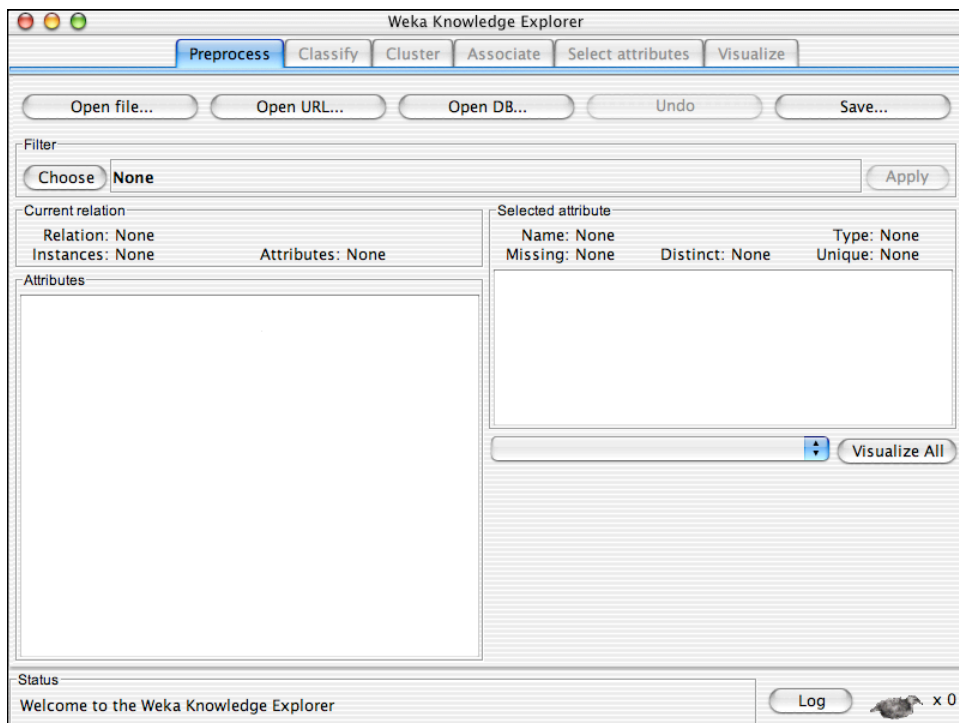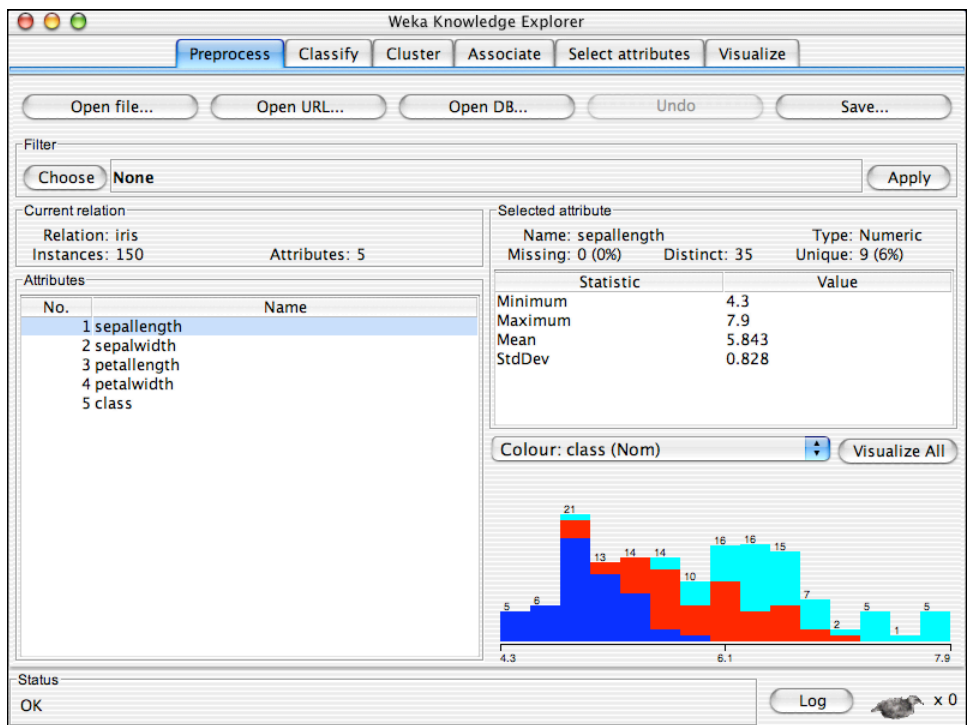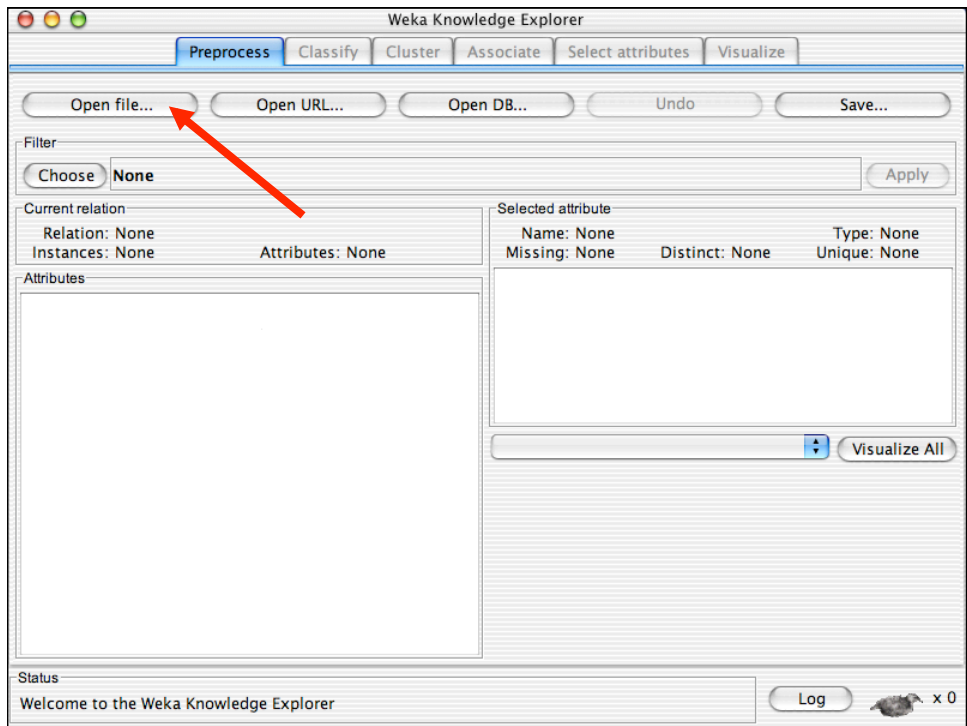
numeric attribute
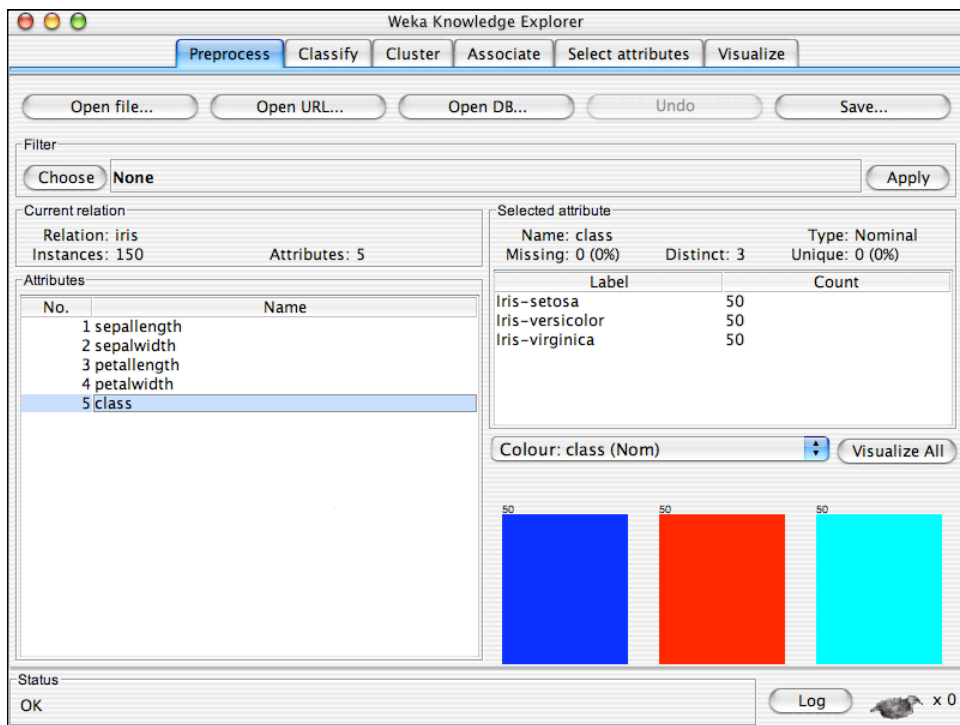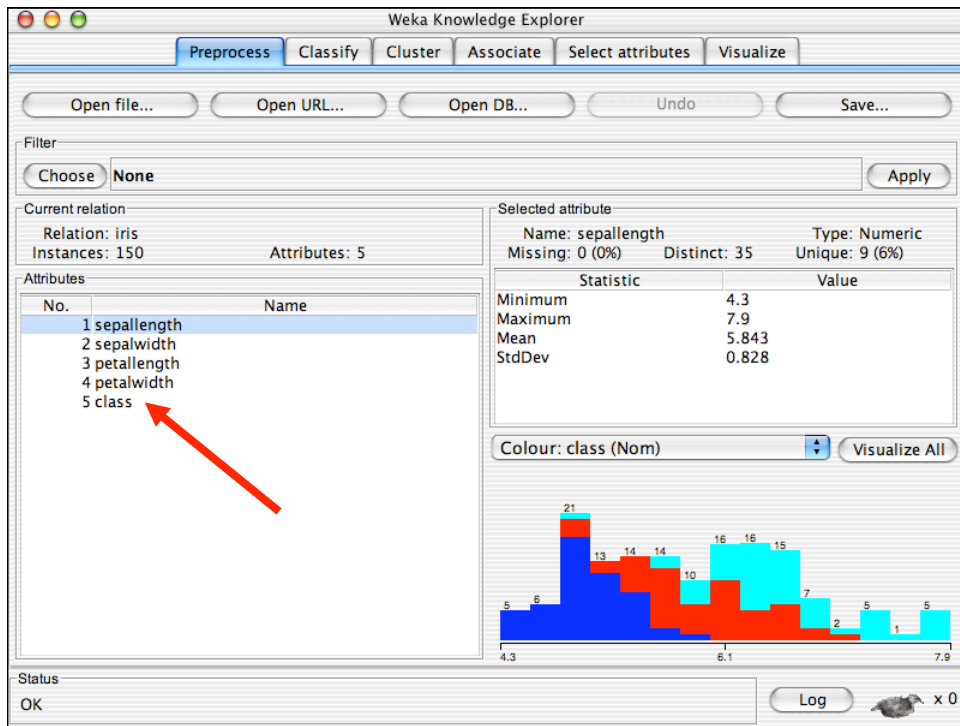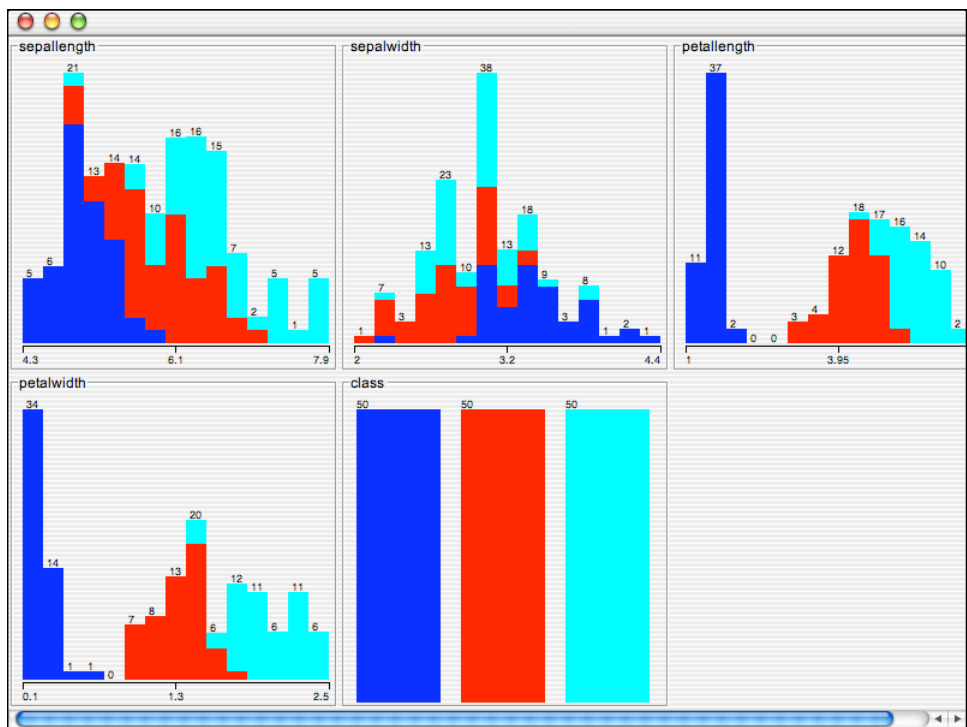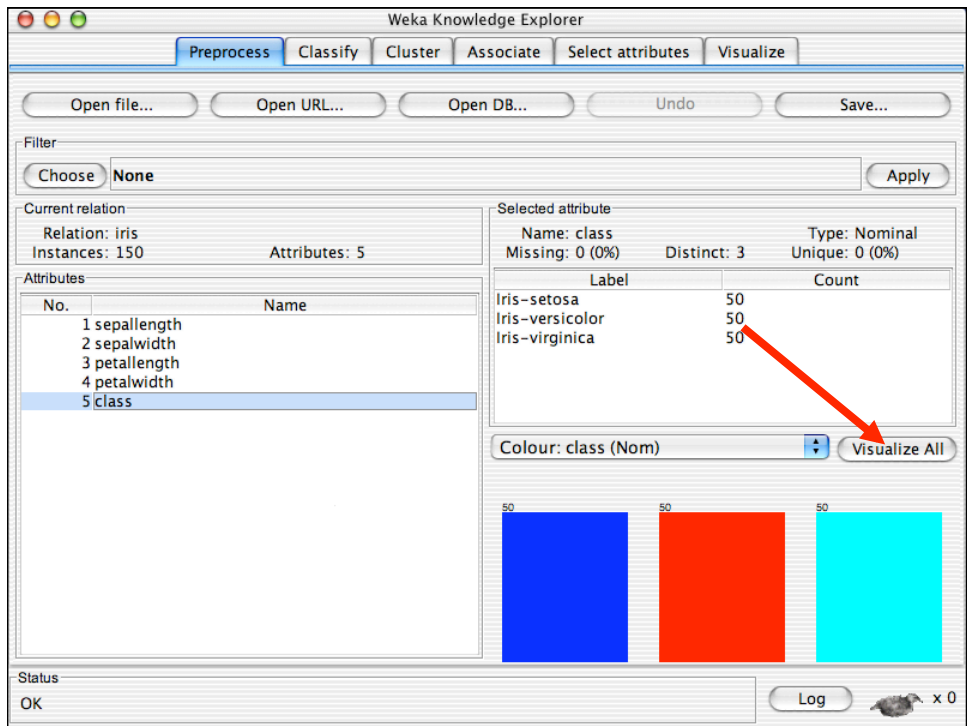
nominal attribute

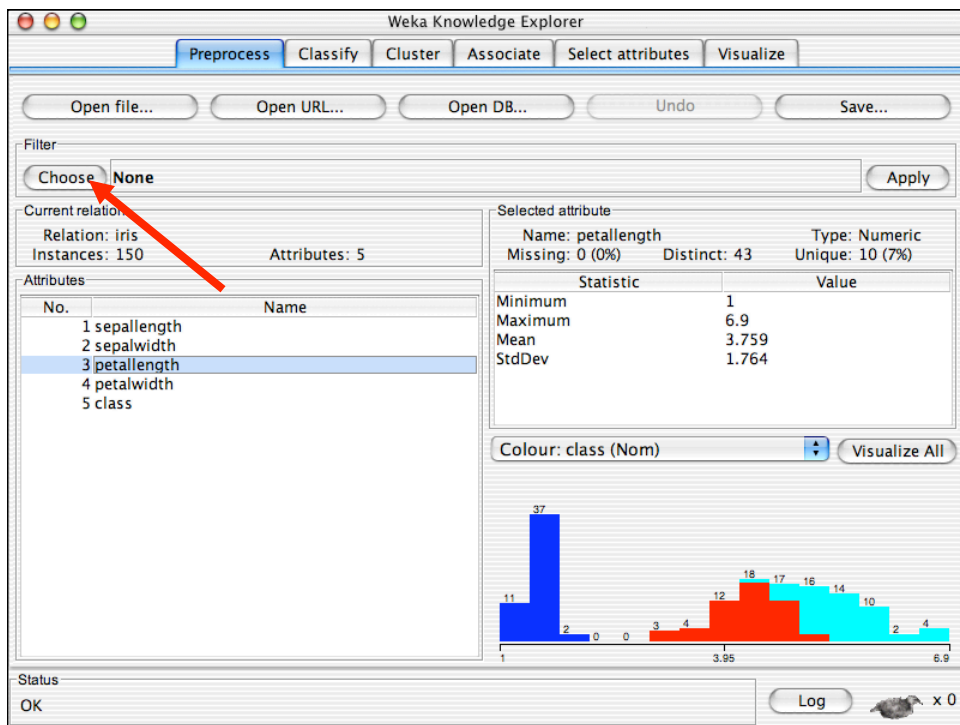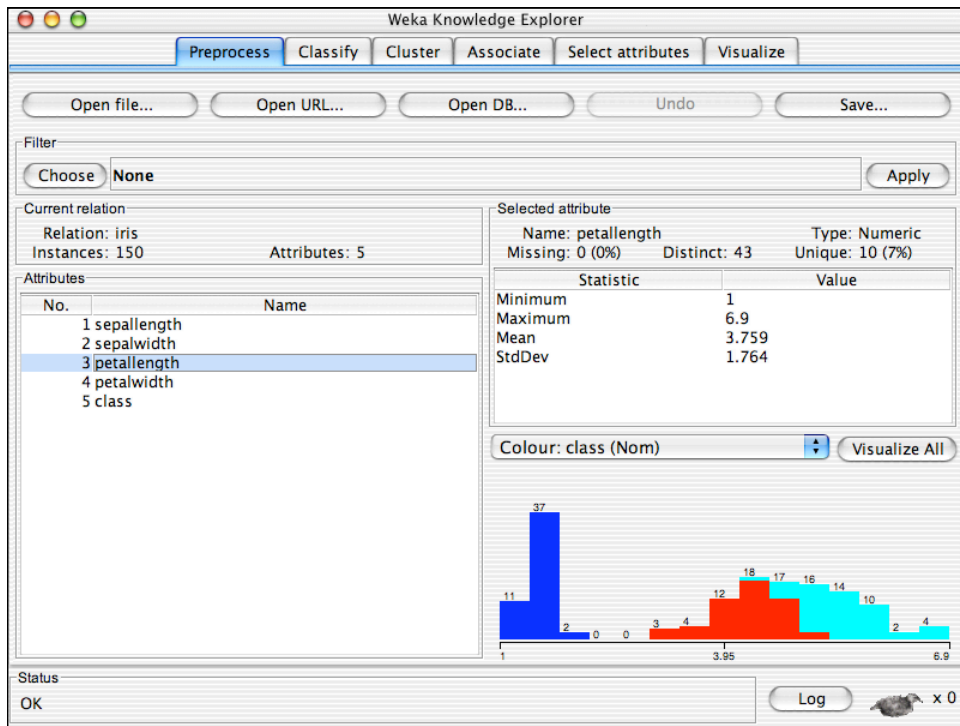Flat file in ARFF format

# Explorer: pre-processing

- **Source**
  - 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**
  - Called "filters"
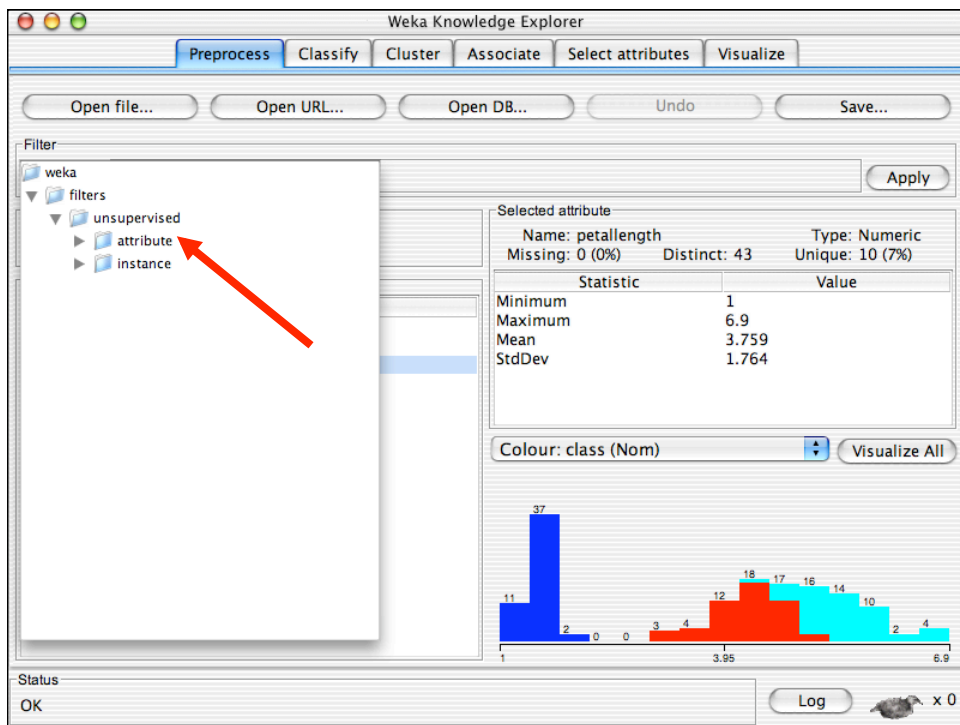  - Discretization, normalization, resampling, attribute selection, transforming and combining attributes, …
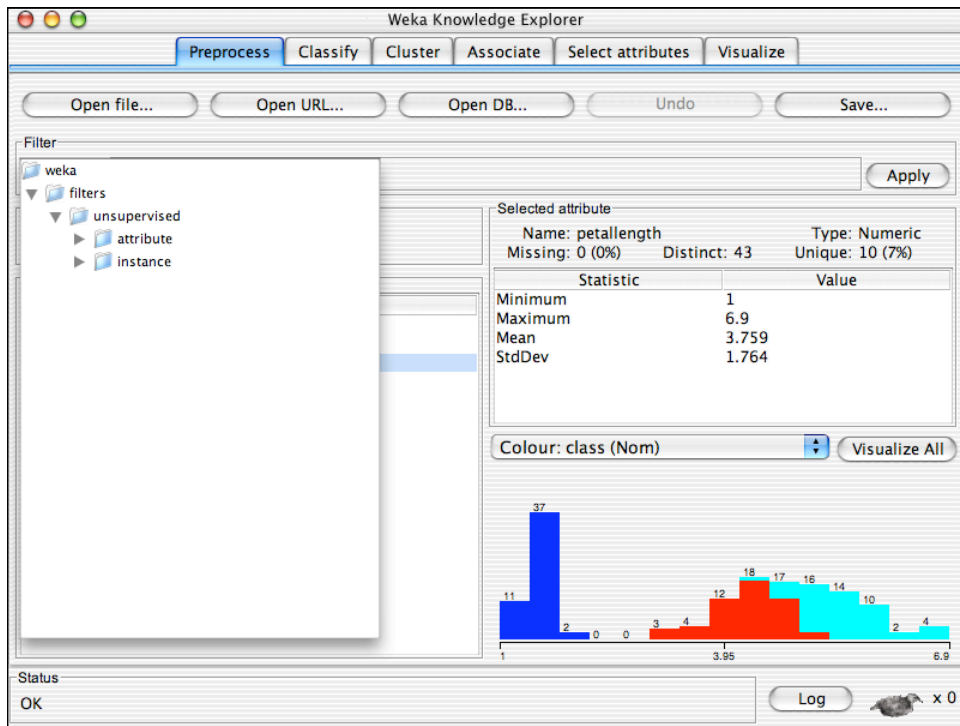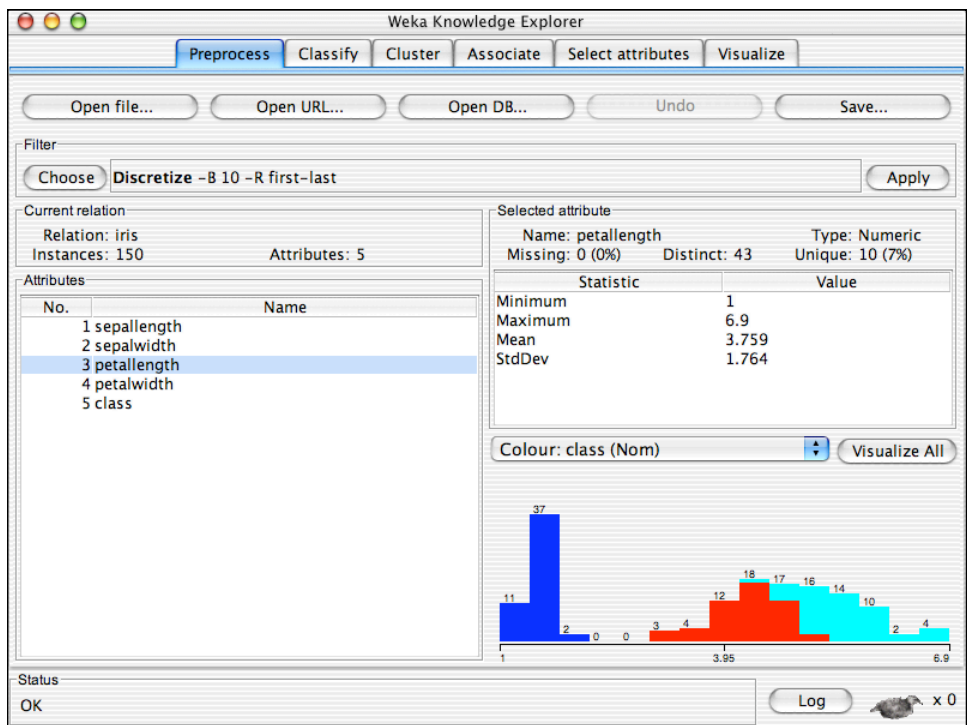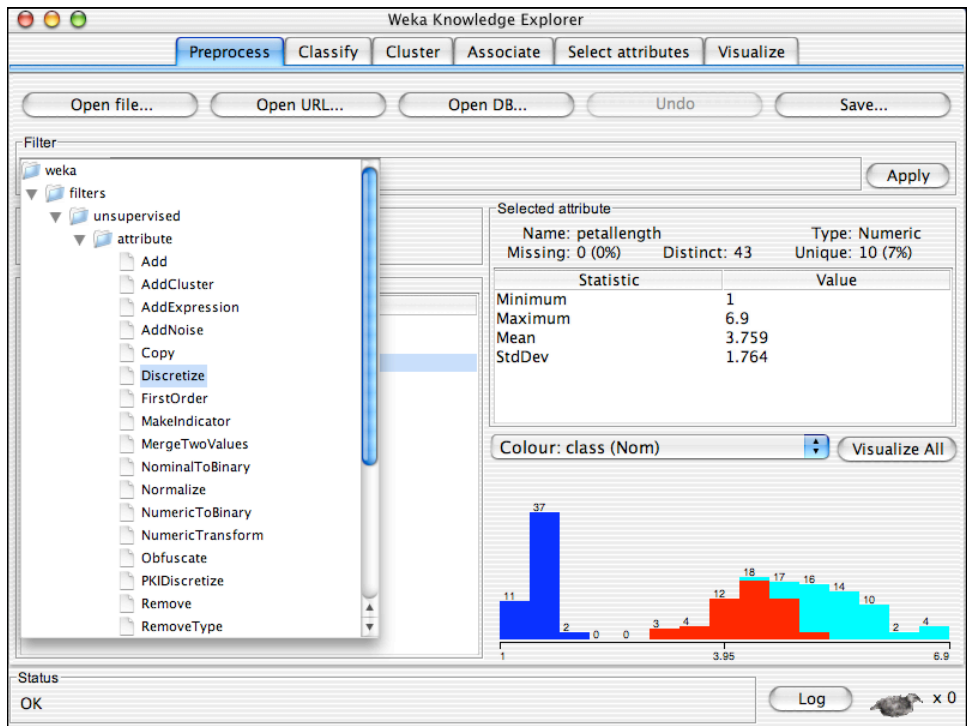
Machine Learning for Data Mining                                                                 5

Machine Learning for Data Mining

# Explorer: 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, …
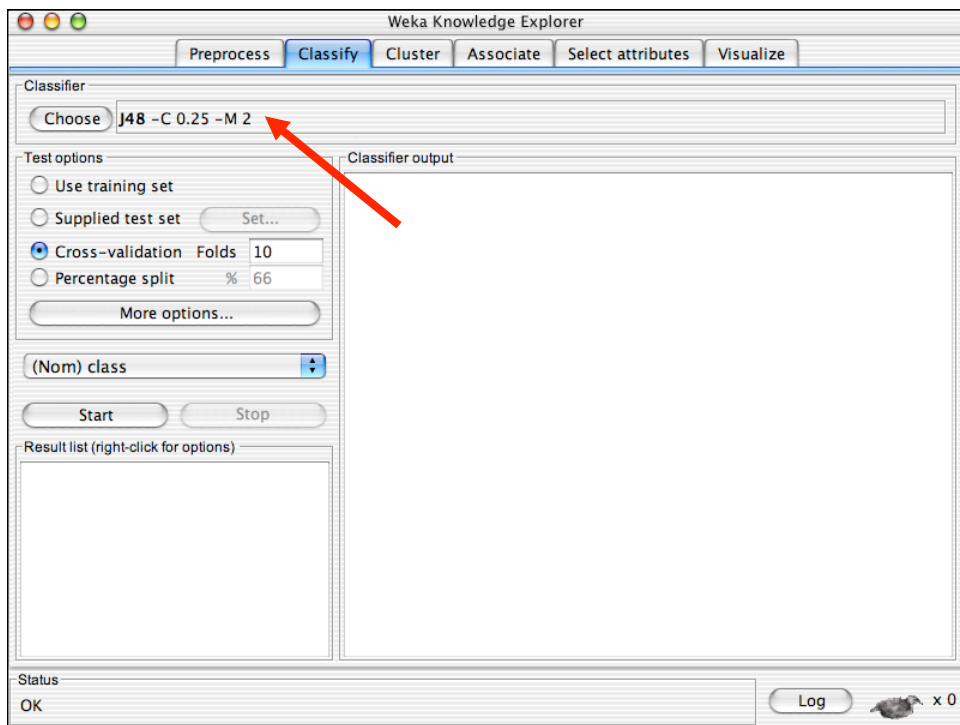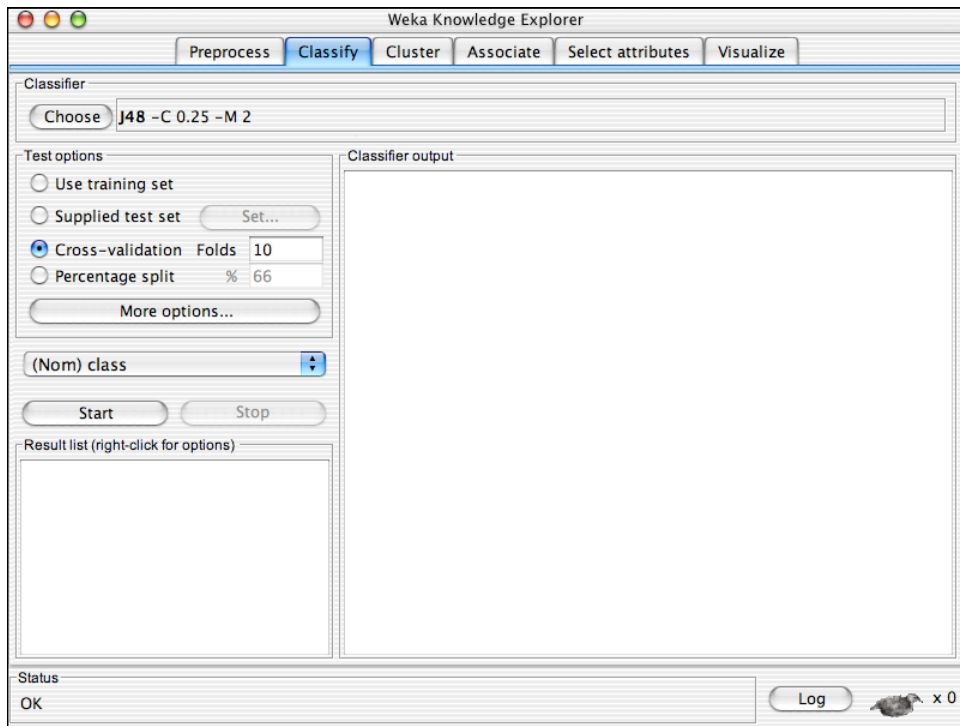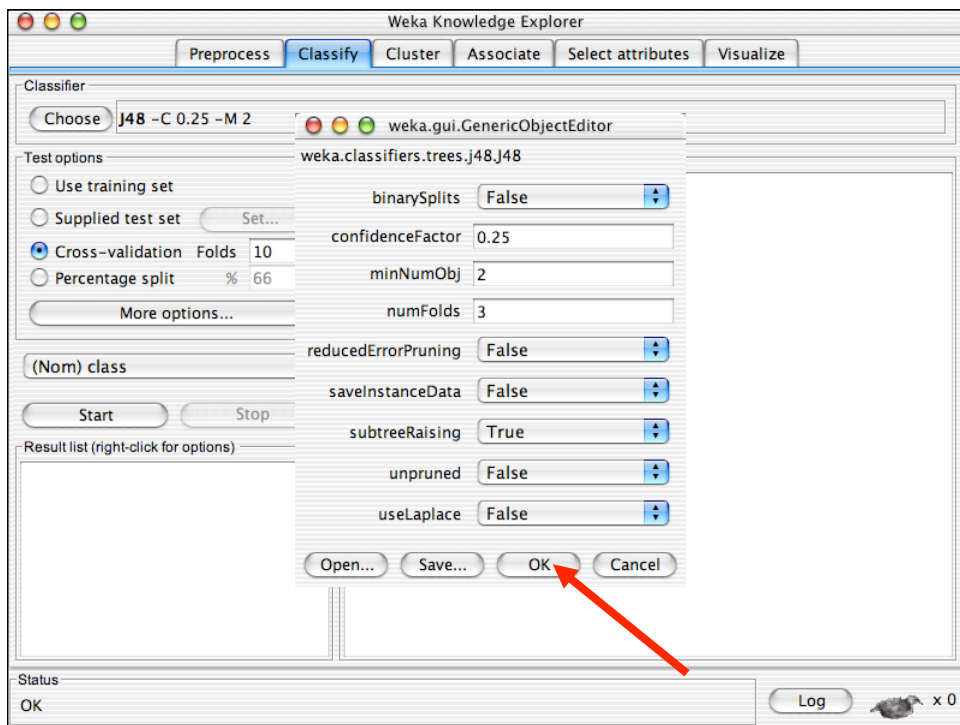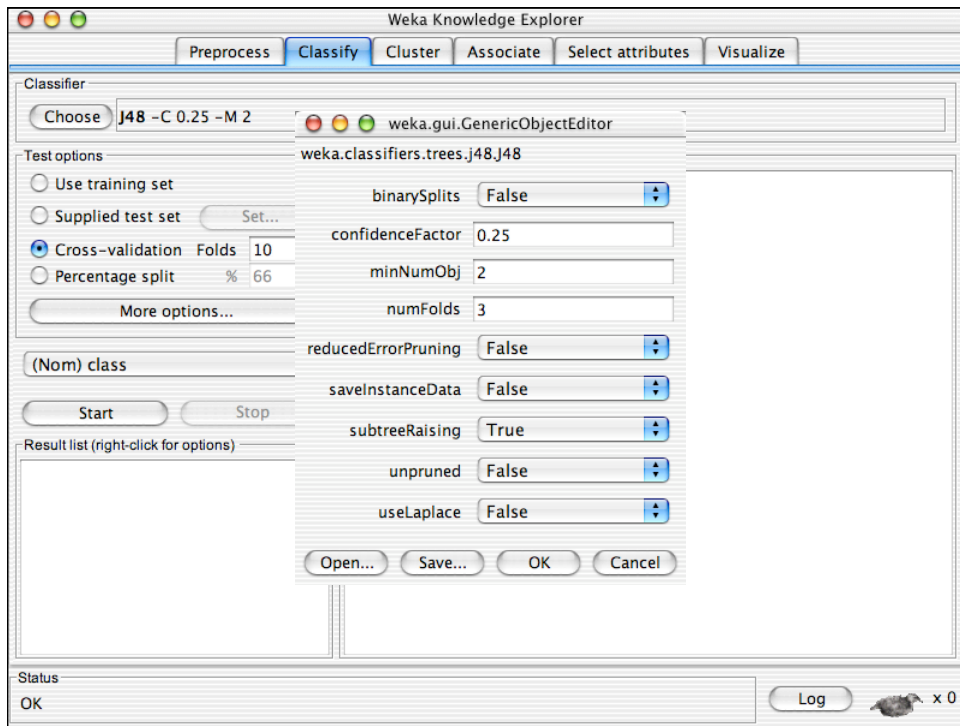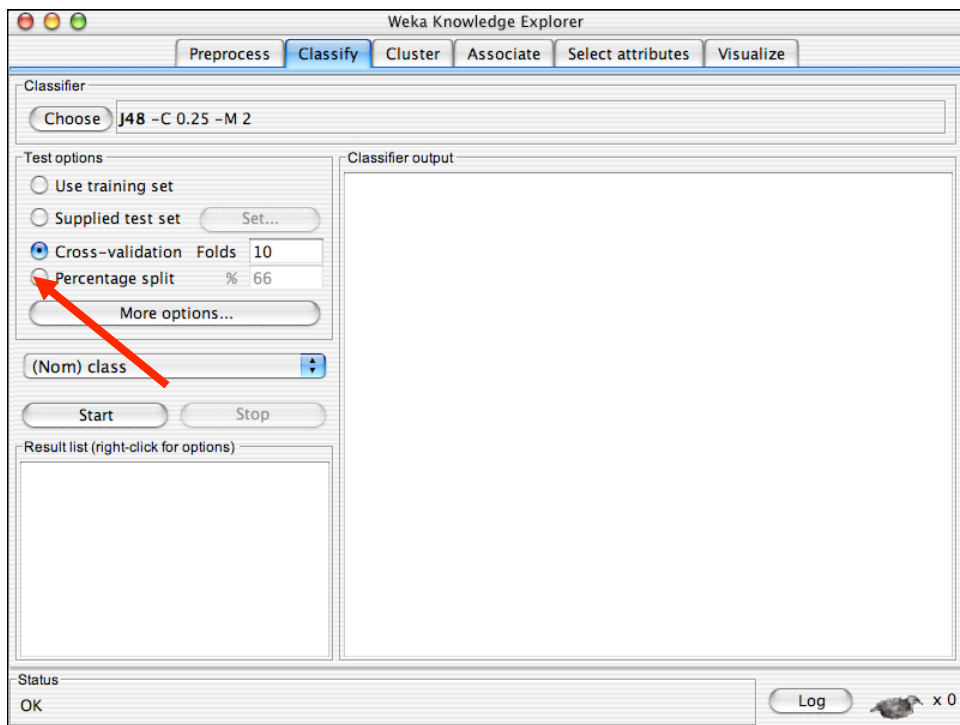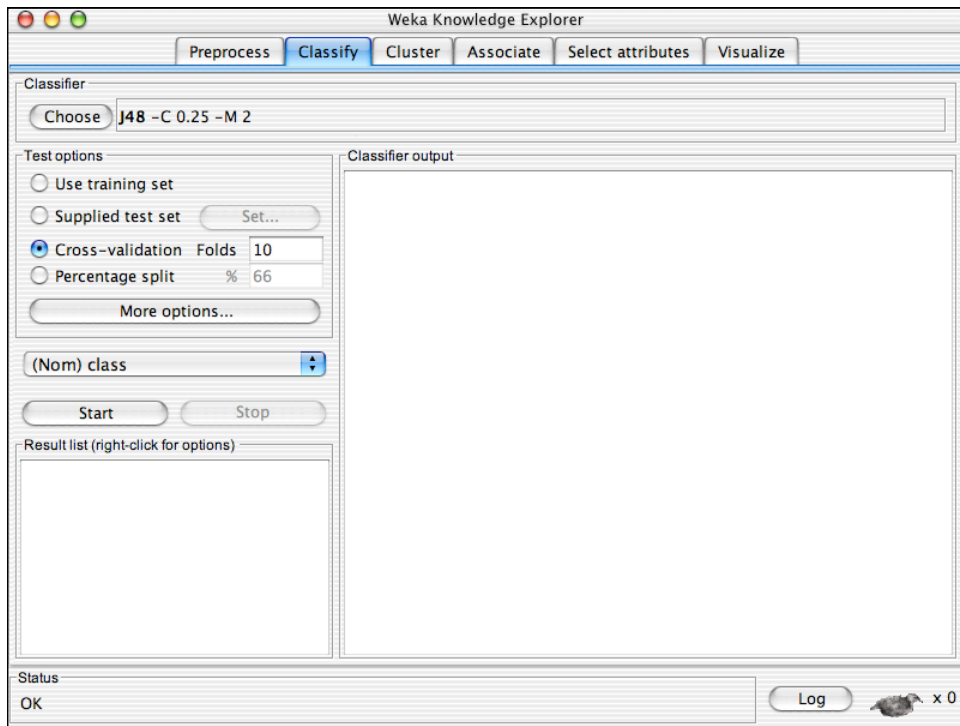
---

**Weka Knowledge Explorer**

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

**Classifier**

Choose | **ZeroR**

**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)**

**Classifier output**

**Status**

OK                          Log      x 0

Machine Learning for Data Mining

Machine Learning for Data Mining

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

consider "TRUE"= iris-virginica and FALSE= everything else

| | Ground truth | |
|---|---|---|
| | FALSE | TRUE |
| detector silent | A =34 | B = 2 |
| detector loud | C= 0 | D = 15 |

| | | | |
|---|---|---|---|
| accuracy | (A+D)/(A+B+C+D) | (34+15)/51 | 96% |
| recall (pd) | D/(B+D) | 15/(2+15) | 88% |
| false alarm (pf) | C/(A+C) | 0/34 | 0% |
| precision | D/(C+D) | 15/(15+0) | 100% |
| f-measure | 2*prec*pd/ (prec+pd) | 2*1*0.88/ (1+0.88) | 94% |

Collect separately for each class.
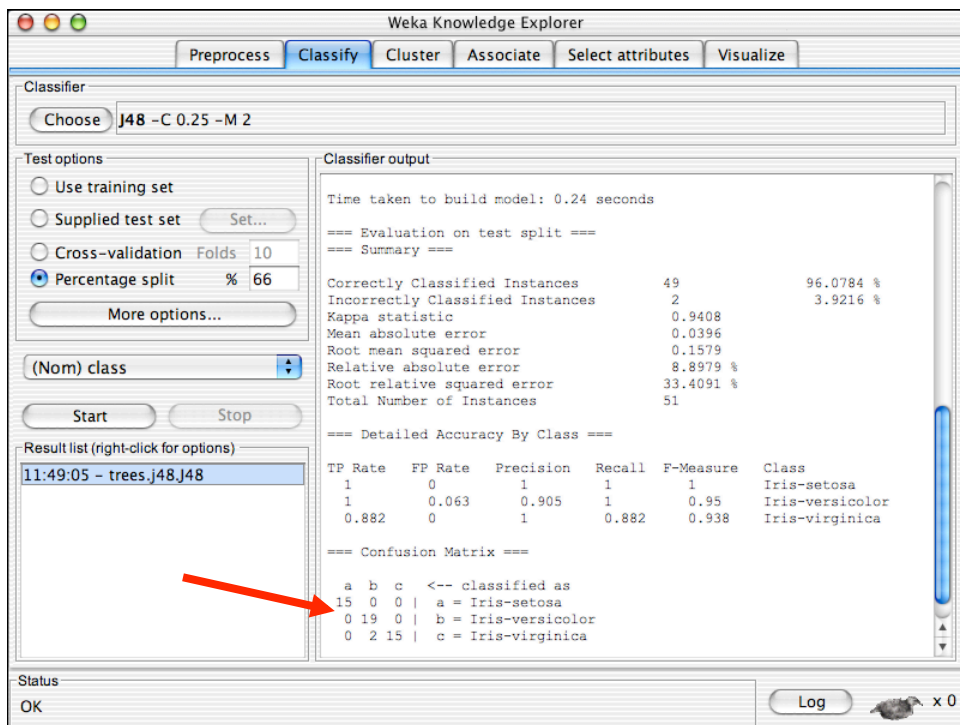Repeat 10 times (re-ordering data) * 10-way
Repeat for each learner * discretizer * x * y * ….

Machine Learning for Data Mining                                                                     30

Machine Learning for Data Mining 34

**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

Result list (right-click for ...)
11:49:05 - trees.j48.J...
14:34:28 - functions...
14:48:05 - bayes.Nai...

Status
OK

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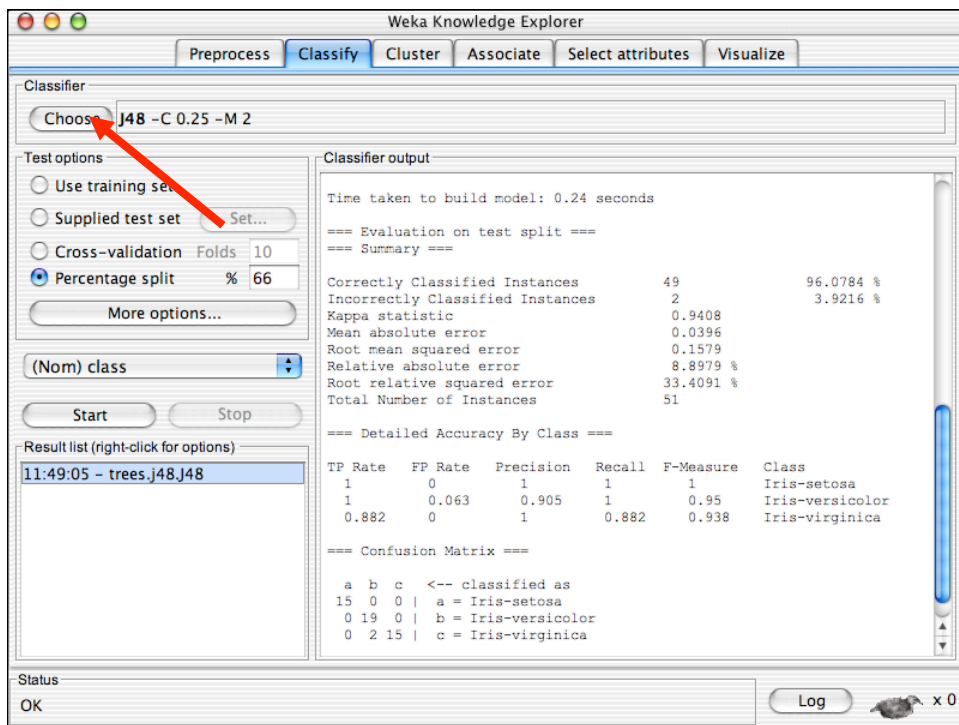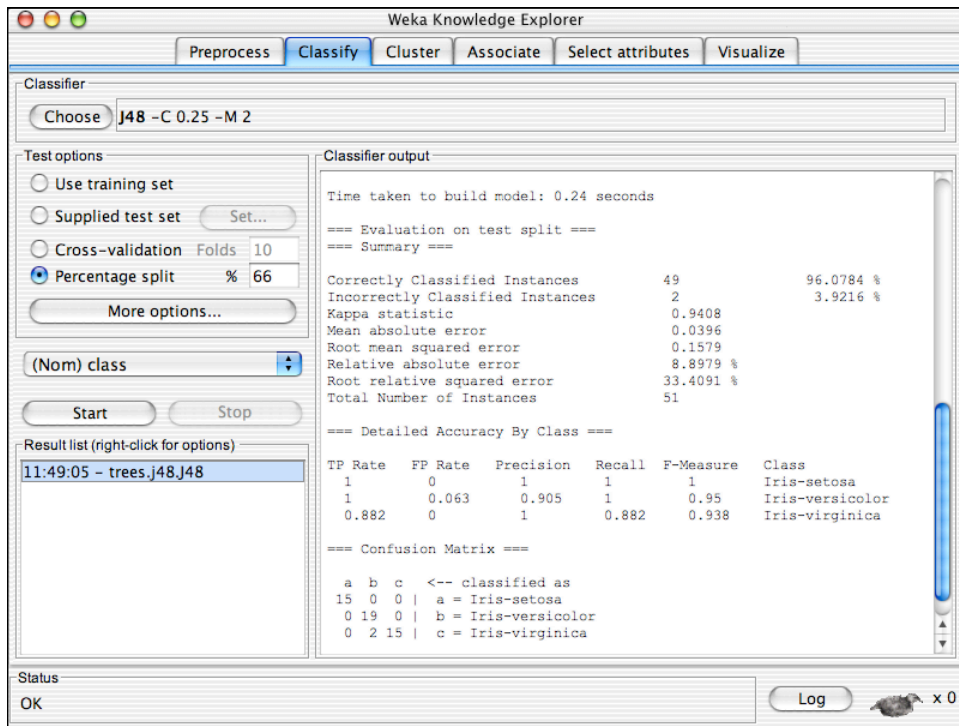
View in main window
View in separate window
Save result buffer

Load model
Save model
Re-evaluate model on current test set

Visualize classifer errors
Visualize tree
Visualize margin curve
Visualize threshold curve  ▶  Iris-setosa
Visualize cost curve           Iris-versicolor
                               Iris-virginica

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

Status
OK

Classifier output
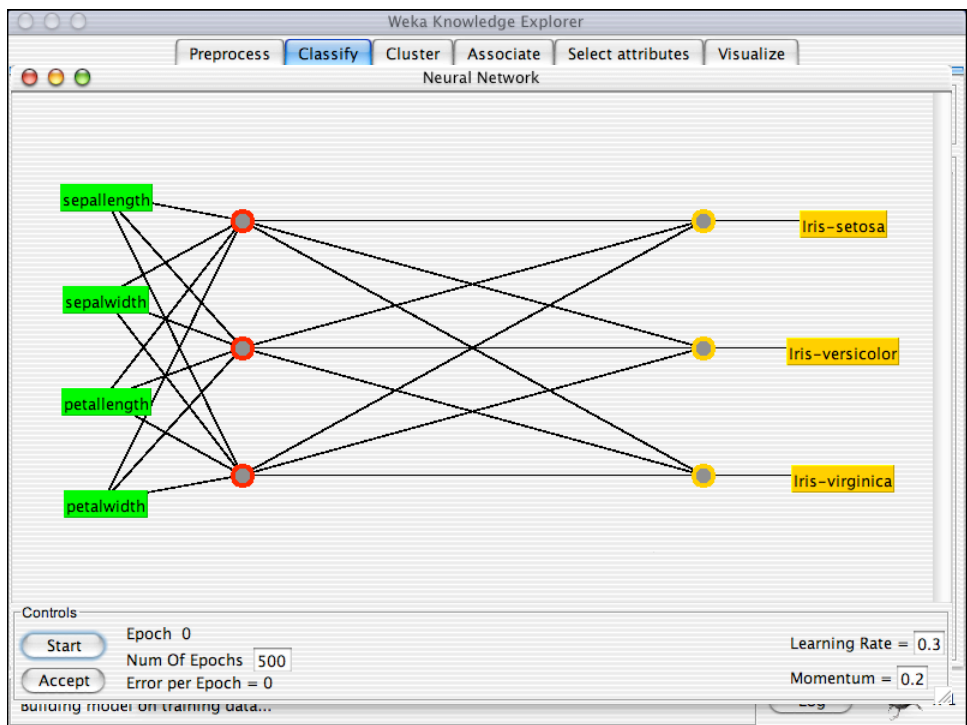
```
=== Evaluation on test split ===
=== Summary ===

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
```
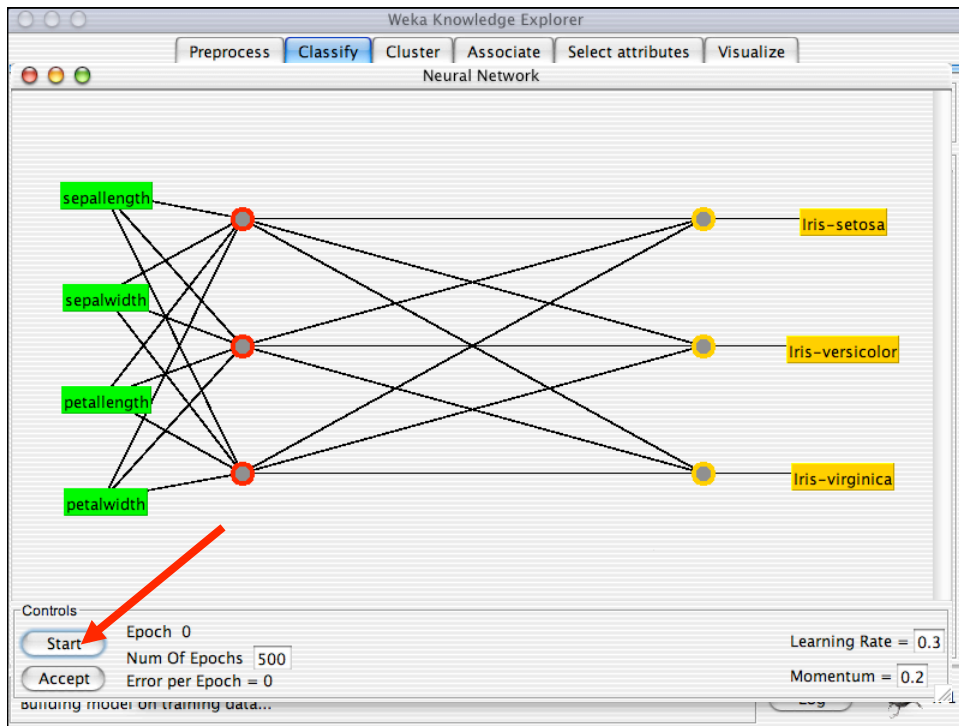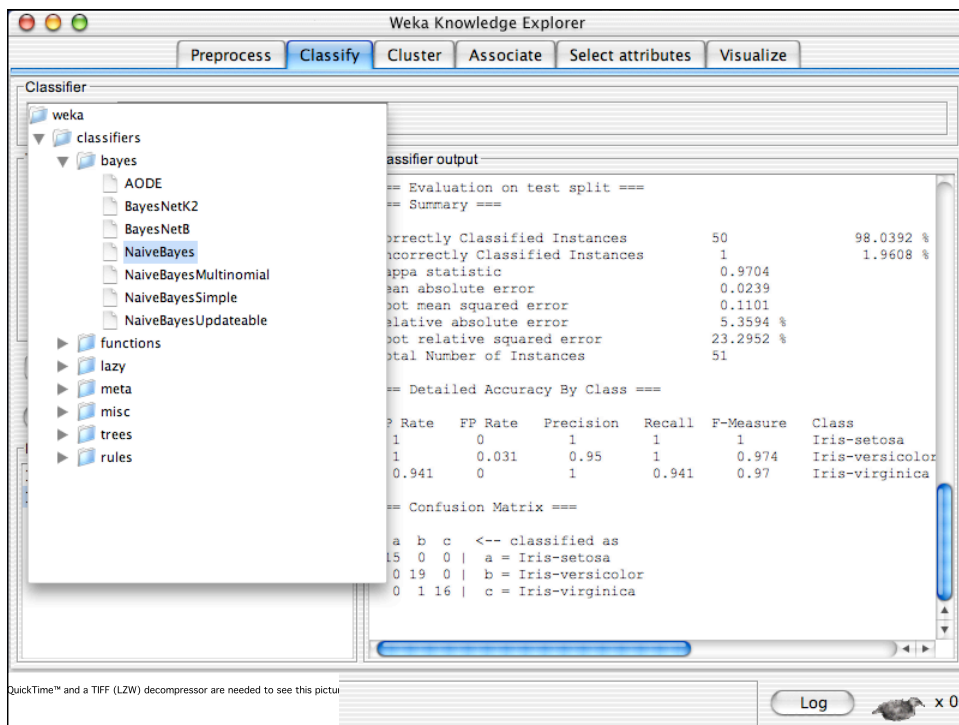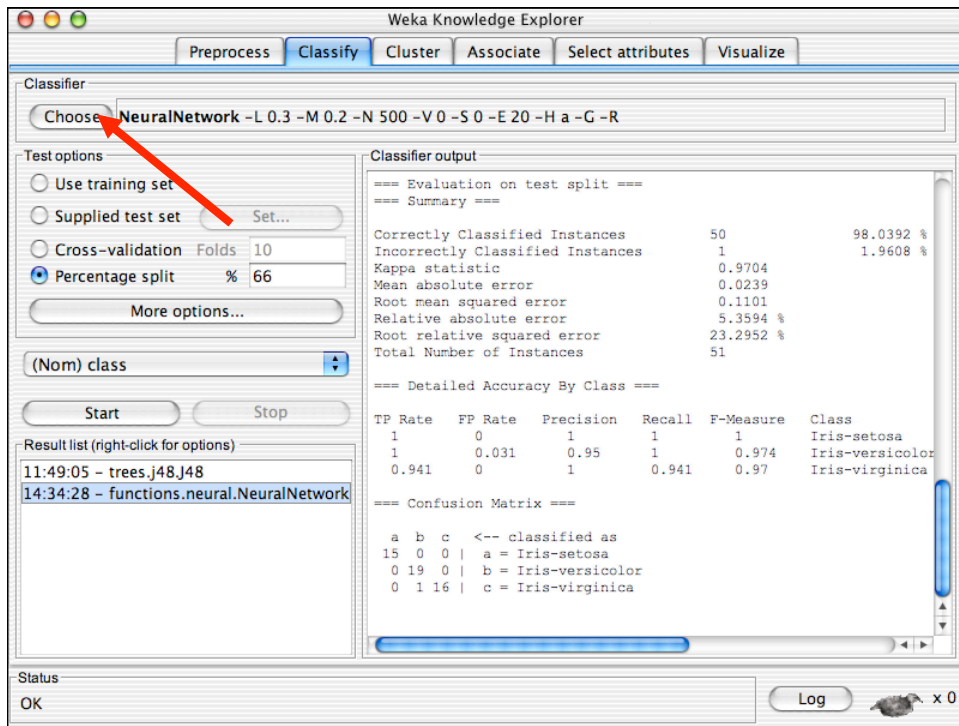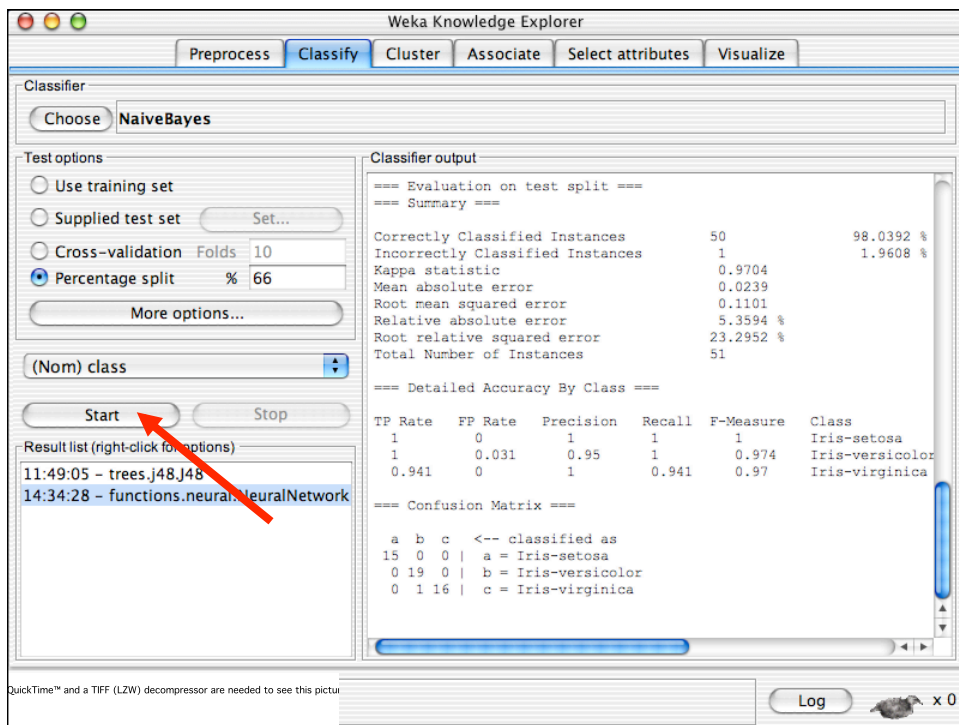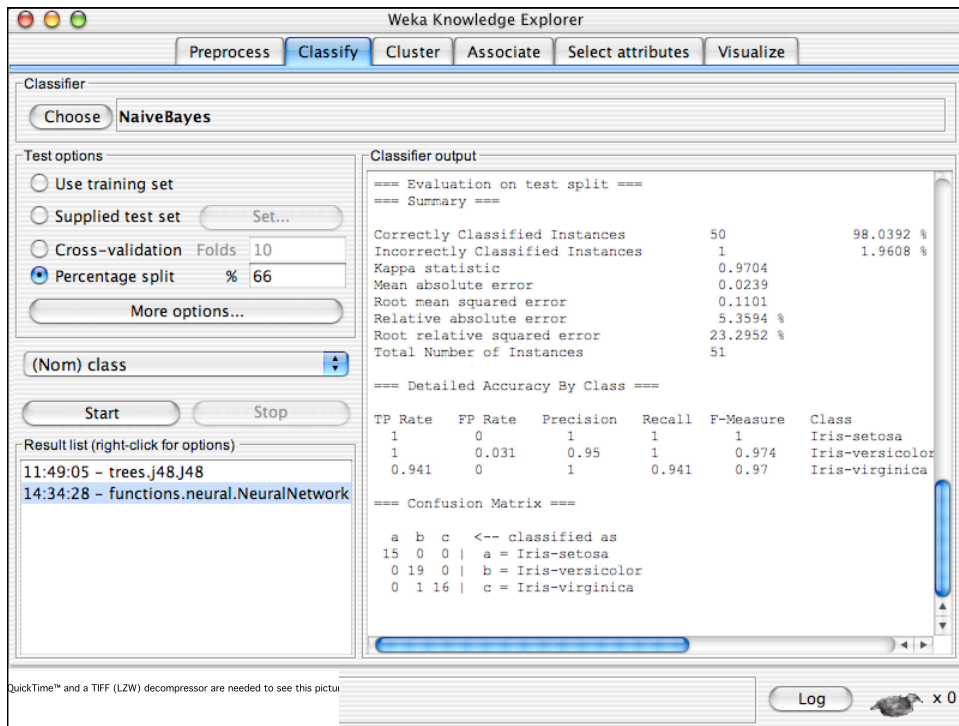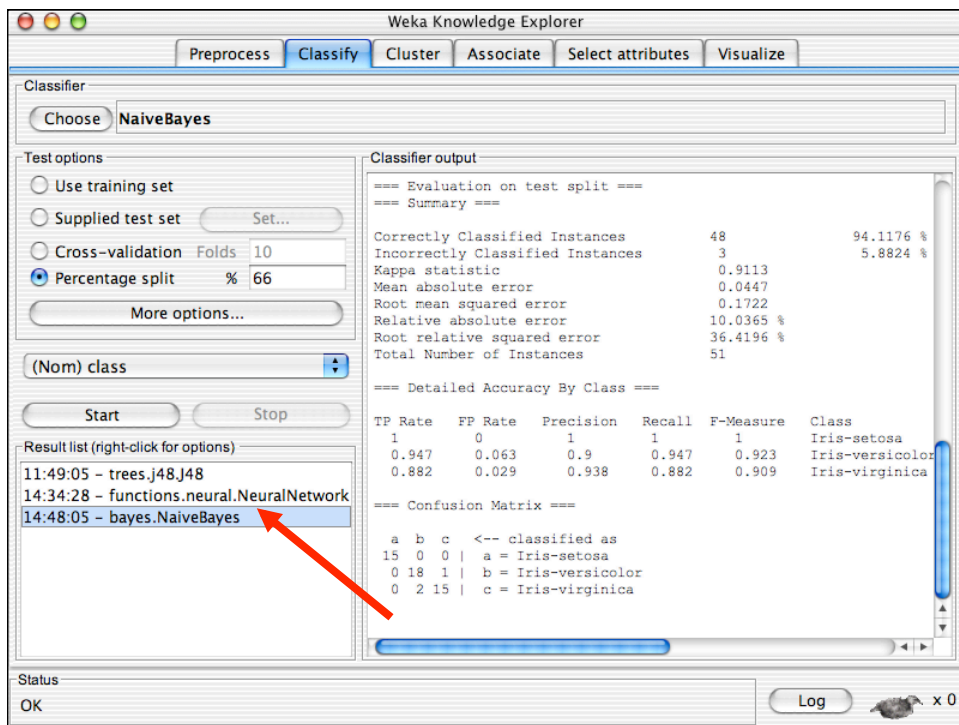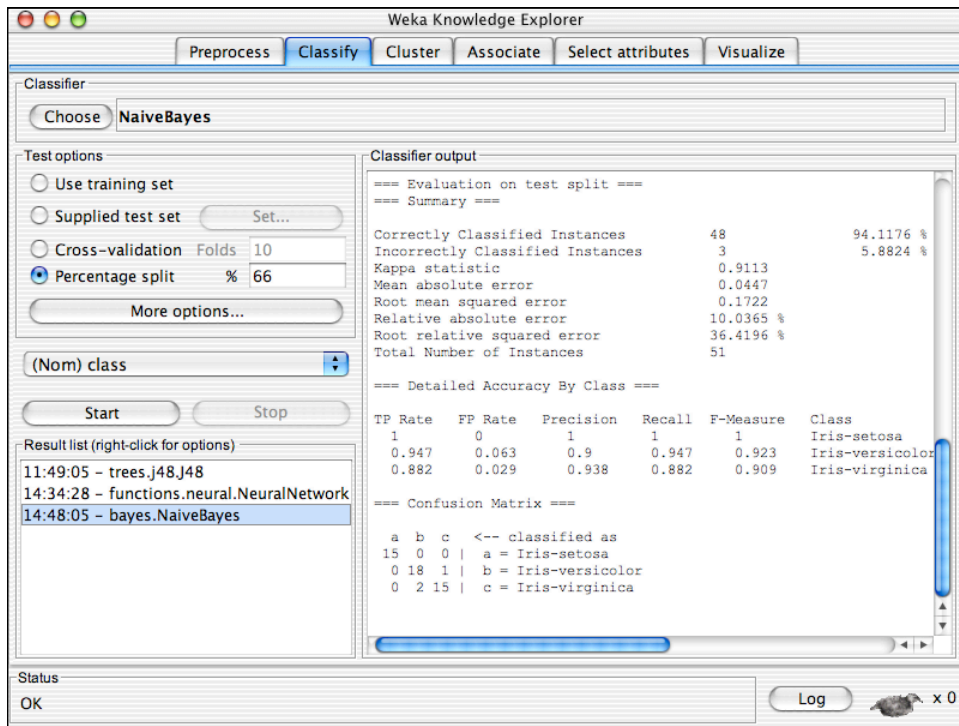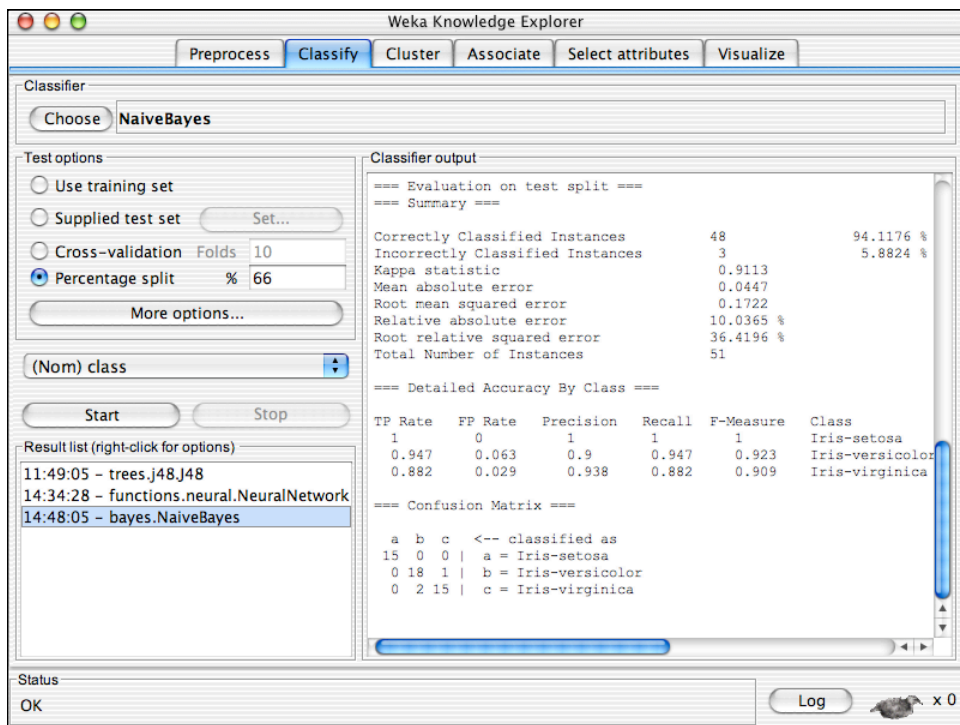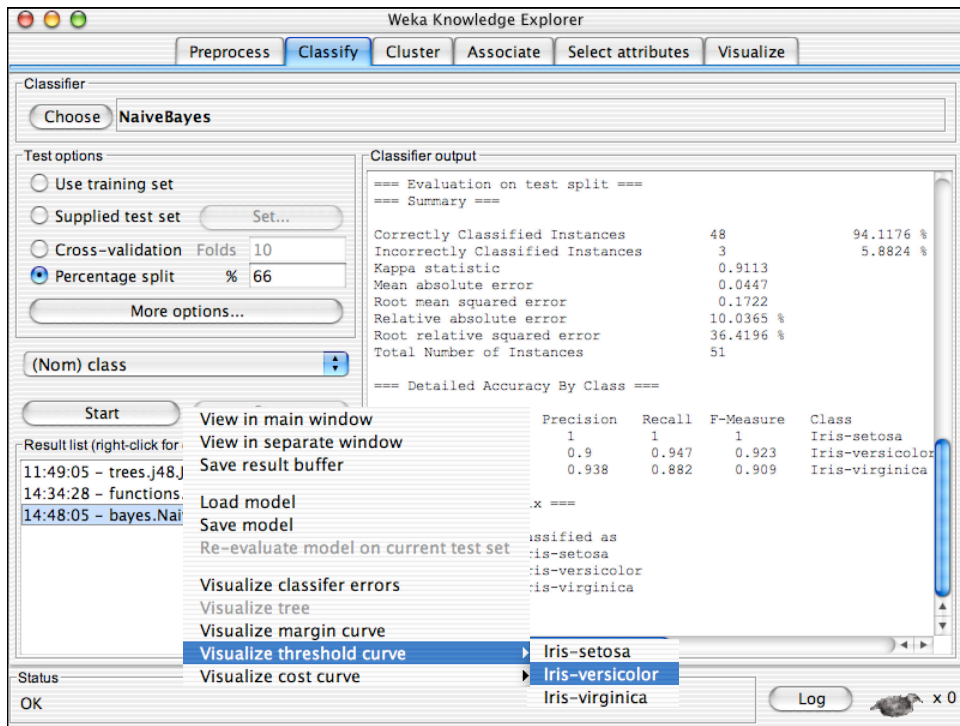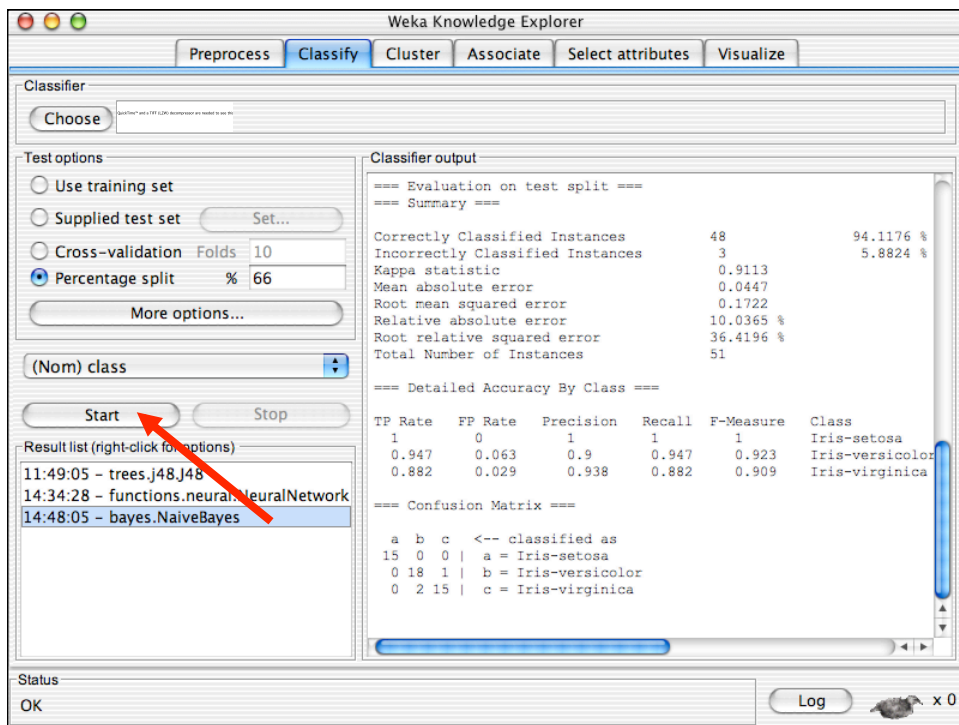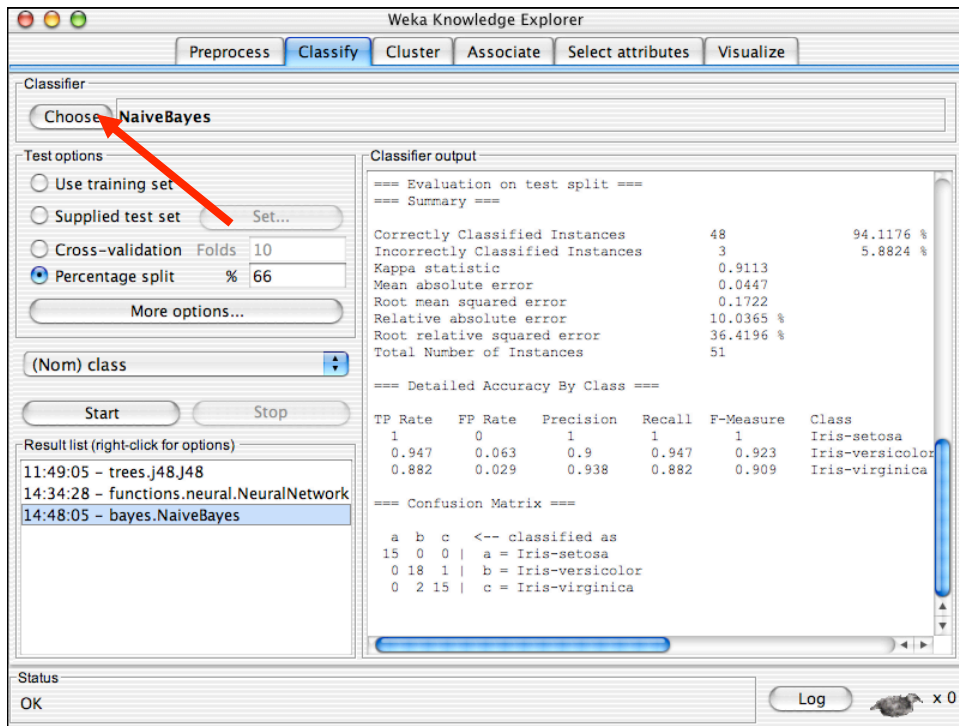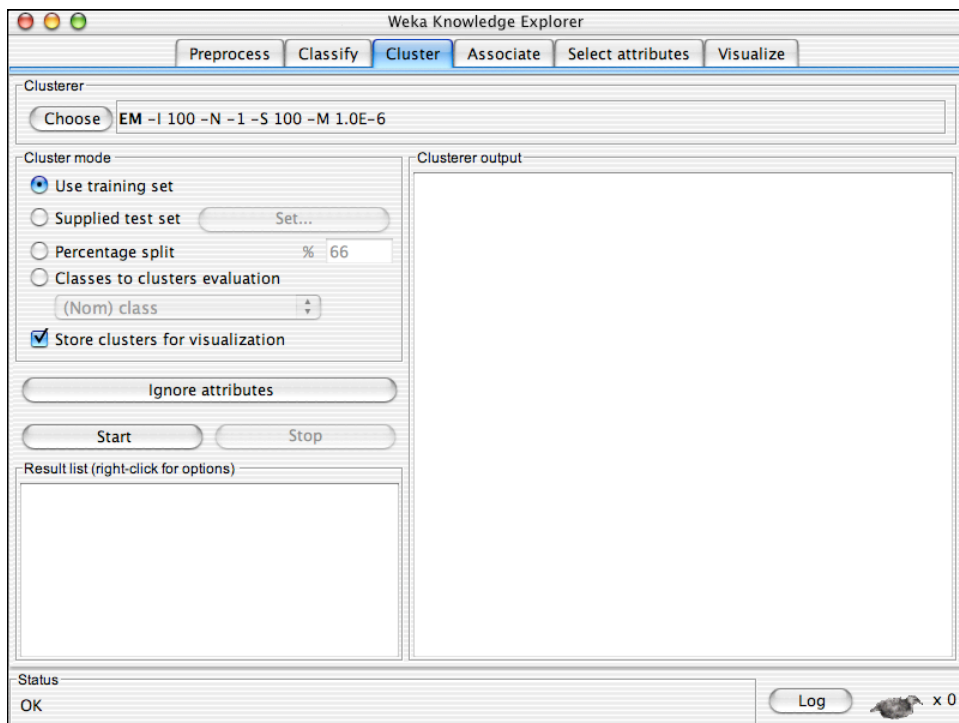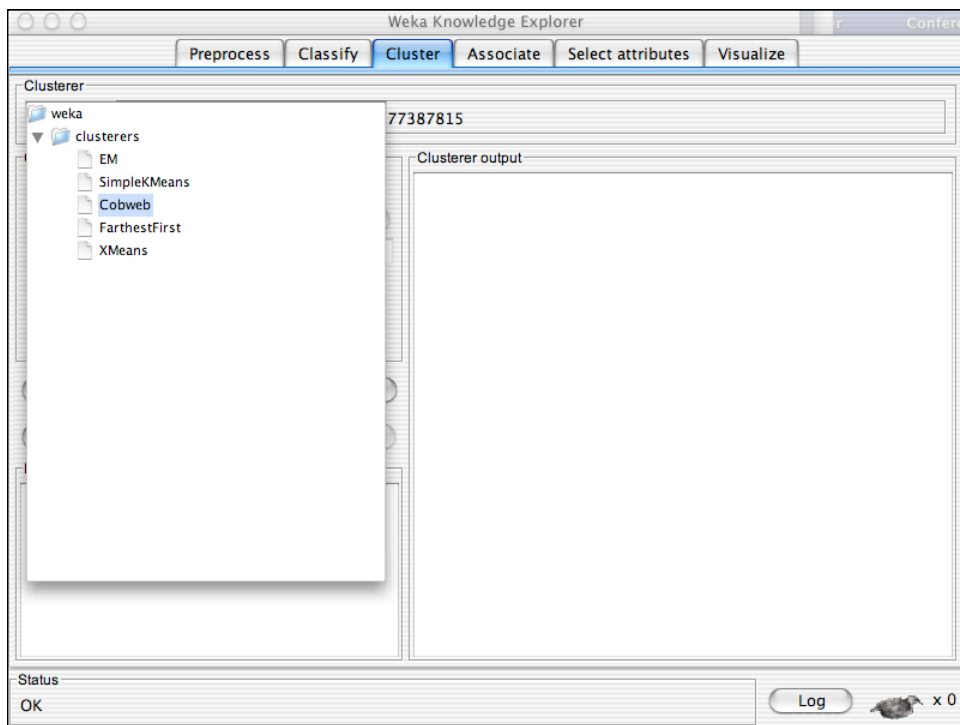
Log    x 0

# Explorer: clustering data

- WEKA contains "clusterers" for finding groups of similar instances in a dataset
- Implemented schemes are:
  - $k$-Means, EM, Cobweb, $X$-means, FarthestFirst
- Clusters can be visualized and compared to "true" clusters (if given)
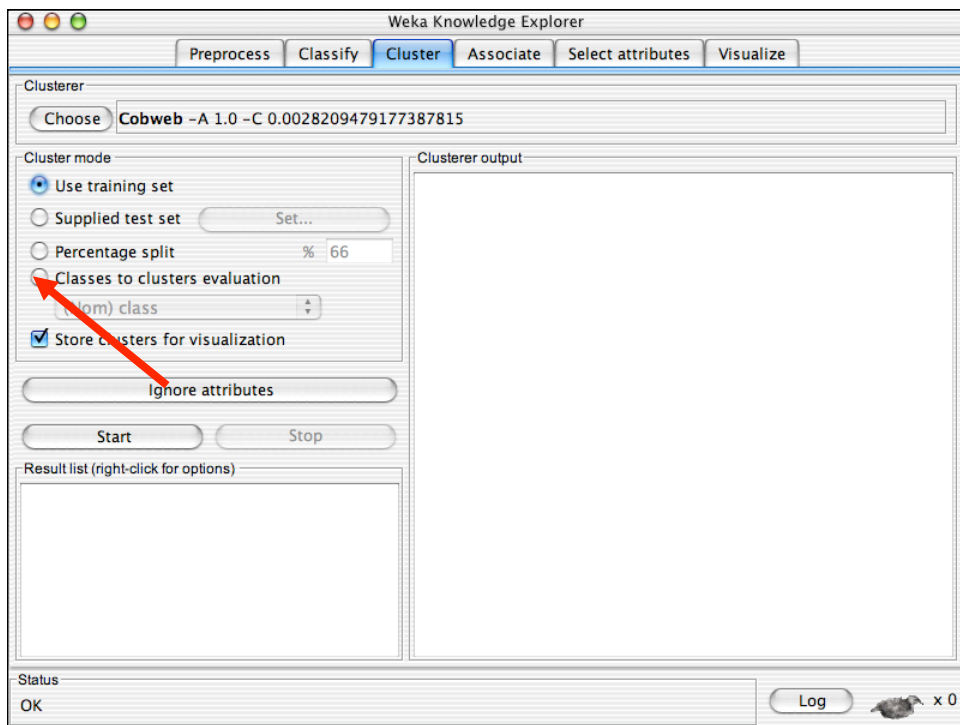- Evaluation based on loglikelihood if clustering scheme produces a probability distribution

Machine Learning for Data Mining

Machine Learning for Data Mining

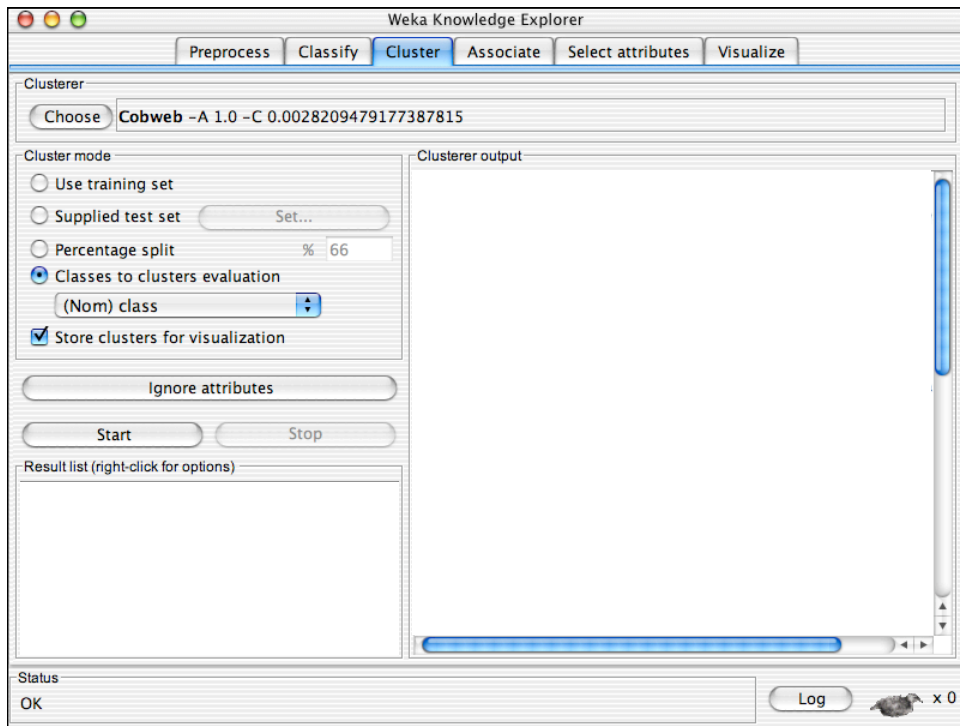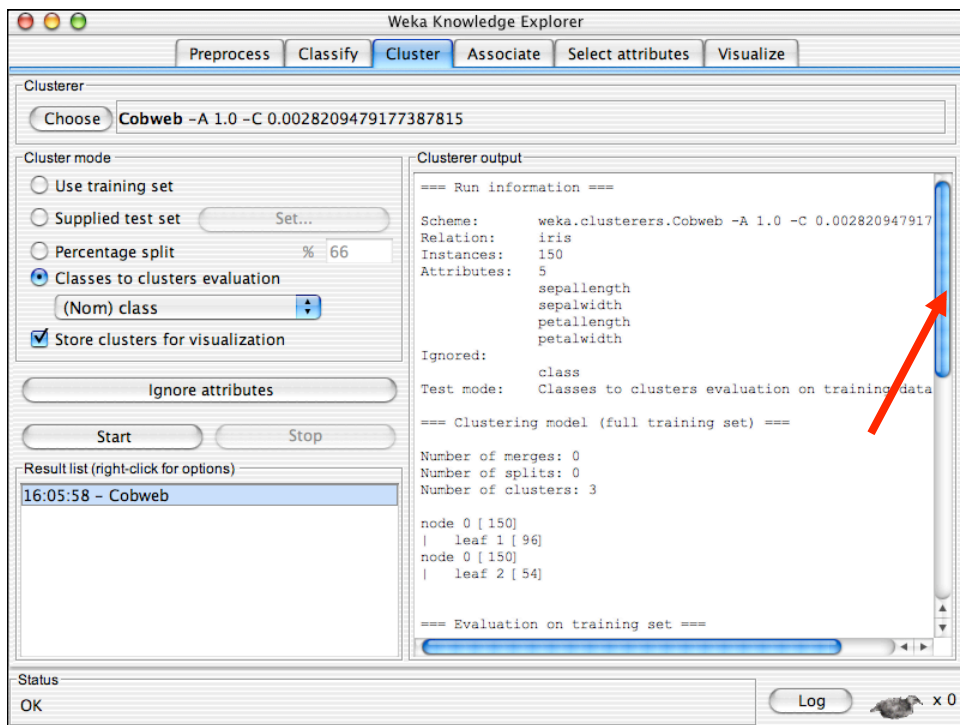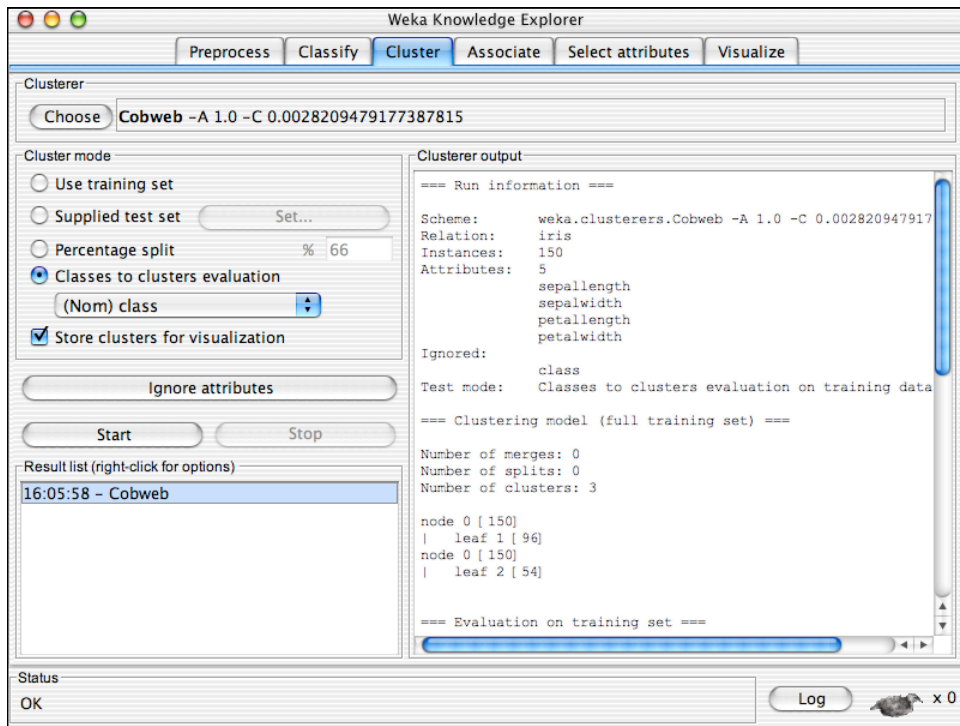# Explorer: finding associations

- 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

Machine Learning for Data Mining

## Explorer: attribute selection

- Panel that can be used to investigate which (subsets of) attributes are the most predictive ones
- Attribute selection methods contain two parts:
  - A search method: best-first, forward selection, random, exhaustive, genetic algorithm, ranking
  - An evaluation method: correlation-based, wrapper, information gain, chi-squared, …
- Very flexible: WEKA allows (almost) arbitrary combinations of these two

Machine Learning for Data Mining                                              49

# Which attribute selector?

- Best: WRAPPER
  - Slow: O(2^N) search through all attribute combinations
  - The "wrapped" learner called to assess each combination
  - Some heuristics to prune the search; but does not scale
- If not WRAPPER
  - Use InfoGain / OneR for very big datasets
  - Use CFS otherwise
- Don't use PCA
  - This is an unsupervised selector
  - So it is uninformed on how dimensions help classification

# Explorer: data visualization

- Visualization very useful in practice: e.g. helps to determine difficulty of the learning problem
- WEKA can visualize single attributes (1-d) and pairs of attributes (2-d)
  - 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

Machine Learning for Data Mining

Machine Learning for Data Mining

Machine Learning for Data Mining 55

Evaluation

# Limitations

- Loads all data into ram prior to learning
  - Problem for large data sets
- Not good for complex experiments
- IMHO, discourages experimentation with new learners
  - The "WEKA effect"
    - Try every learner till something works
- Still, very useful for
  - Initial investigations
  - Learning data mining
  - Or as a sub-routine of other systems

# Alternate tools: Orange



Written in Python

Simpler specification (but see WEKA's KnowledgeFlow Environment).

Also, less community support/ debugging. So sometimes frustrated by random bugs

# Alternate tools: RapidMiner



Experiments specified in an XML tree syntax

In theory, possible to share experimental descriptions

# Alternate tools: OurMine

```
Java=$Base/lib/java
Weka="java -Xmx2048M -cp $Java/weka.jar "
Clusterers="java -Xmx1024M -jar $Java/Clusterers.jar "
Reducers="java -Xmx1024M -jar $Java/Reduce.jar "

nb() {
    local learner=weka.classifiers.bayes.NaiveBayes
    $Weka $learner -p 0 -t $1 -T $2
}

nb10() {
    local learner=weka.classifiers.bayes.NaiveBayes
    $Weka $learner -i -t $1
}

j48() {
    local learner=weka.classifiers.trees.J48
    $Weka $learner -p 0 -C 0.25 -M 2 -t $1 -T $2
}
```
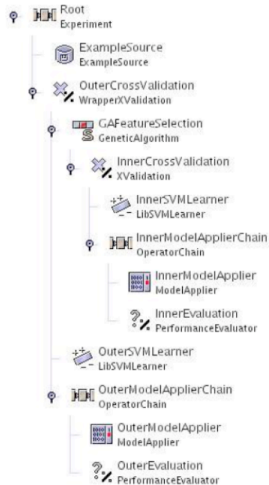
Forget the visuals.

Make WEKA a sub-routine inside Bash script

Now you can mix WEKA's JAVA with learners written in your favorite language.

But how do you find the magic command strings?

# Why go to all that trouble?

```
analysis1(){
   local origdata=$1
   local outstats=$2
   local nattrs="2 4 6 8 10 12 14 16 18 20"
   local learners="nb10 j4810 zeror10 oner10 adtree10"
   local reducers="infogain chisquared oneR"
   local tmpred=$Tmp/red
   echo "n,reducer,learner,accuracy" > $outstats

   for n in $nattrs; do
      for reducer in $reducers; do
         $reducer $origdata $n $tmpred
         for learner in $learners; do
            accur=`$learner $tmpred.arff | acc
            out="$n,$reducer,$learner,$accur"
            blabln $out
            echo $out >> $outstats
         done
      done
   done
}
```

Complex experiments, specified succinctly.

Experiments can now be reviewed, audited, by others.

Also, in 12 months time when Reviewer2 wants a tiny extension to the old paper, you don't have to remember all that clicking you did: just rerun the script.