# Machine Learning the WEKA directions
## Week-7

# WEKA IS

- ✓ Machine Learning toolkit

- ✓ Open Source (Java)

- ✓ Well documented + huge community

- ✓ Provides API, command line and GUI-swing tools

- ✓ Relatively easy to learn

- ✓ Runnable on a remote server so you can "dumb terminal" your laptop and keep your data in one place!

# WEKA is not...

- A complete replacement for R/Matlab

- Optimized out of the box for multiple CPUs / compute farms

- An excuse to ignore the method details of a clustering/classification algorithm
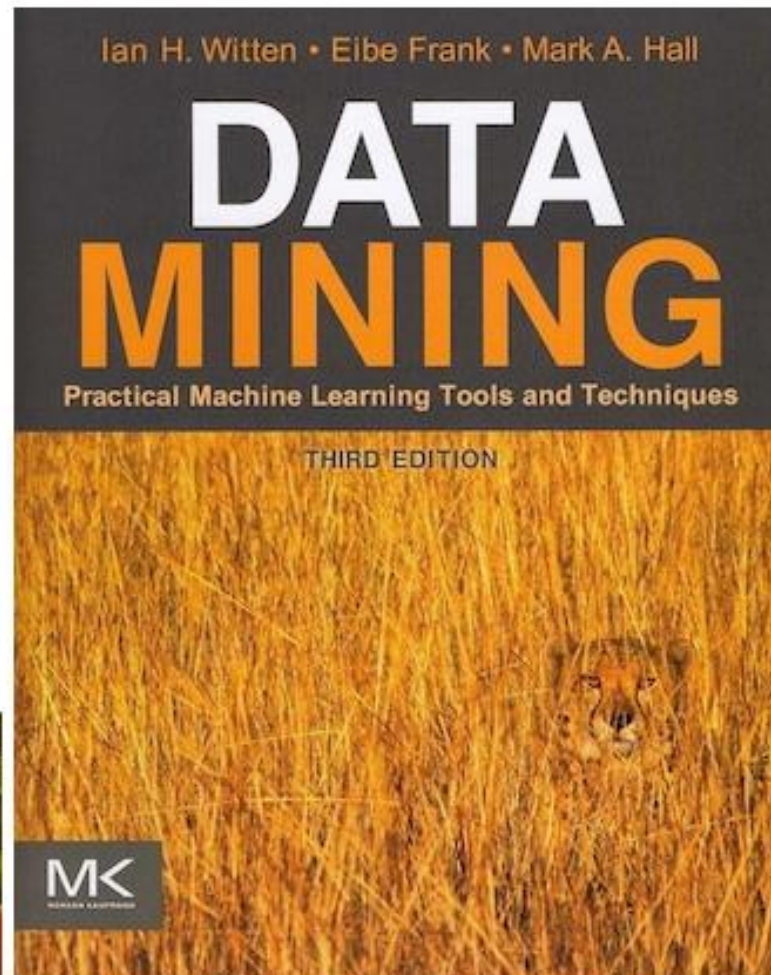
# Data Mining: Practical Machine Learning Tools and Techniques (Third Edition)

**Ian H. Witten**, **Eibe Frank**, **Mark A. Hall**

Morgan Kaufmann
January 2011
629 pages
Paper
ISBN
978-0-12-374856-0



Eibe Frank and Ian Witten



Click here to order from Amazon.com

# Personal WEKA dataset

- ✓ 21k+ variables
- ✓ 14k+ subjects phenotyped
- ✓ 9k+ subjects genotyped 500k Affymetrix
- ✓ 54M recorded phenotype values of widely varying types

- ✓ Even the simplest correlation matrix
  - ✓ 20k * 20k = **400M** comparisons *before* including SNPs

# WEKA basics
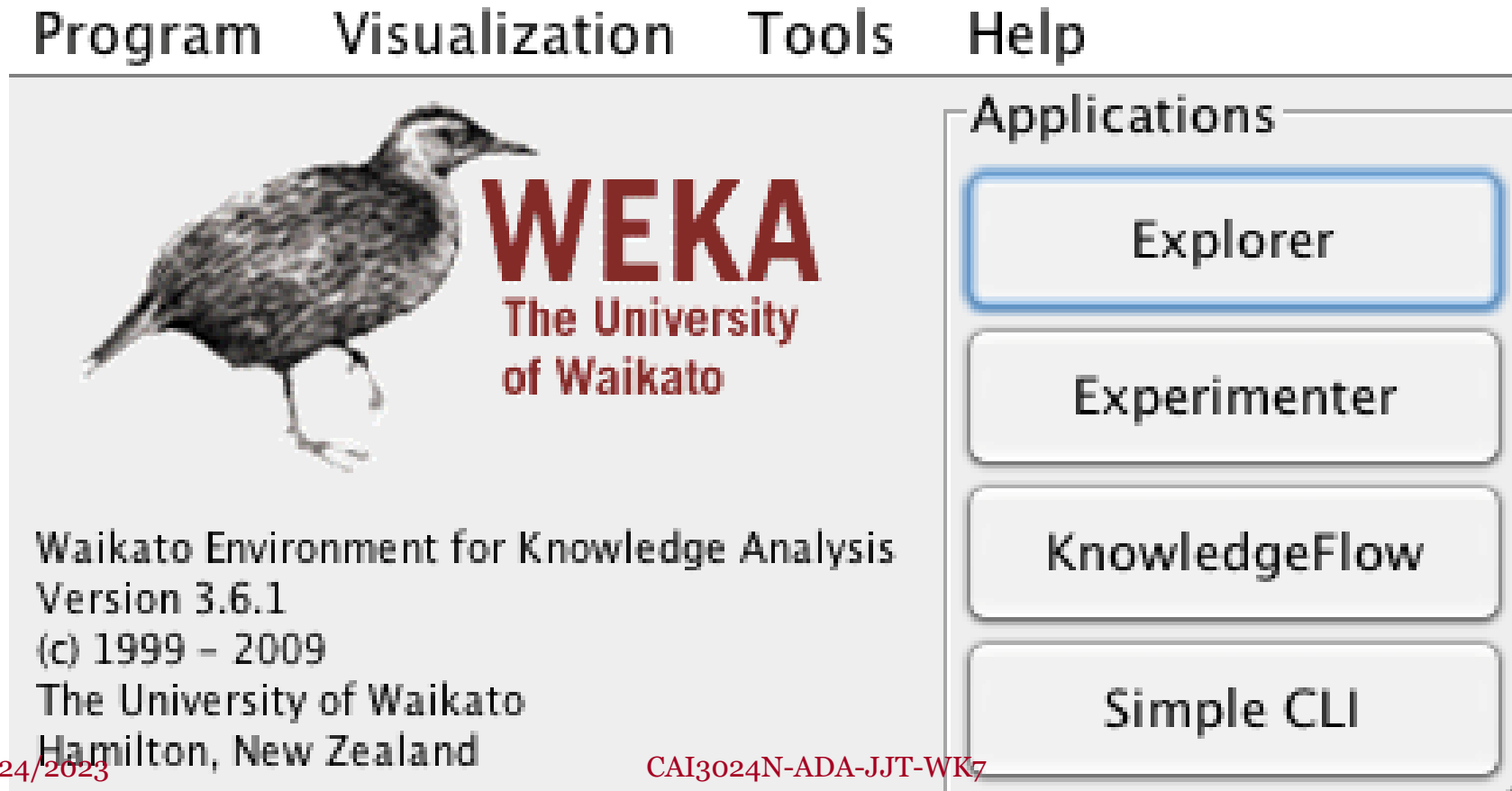## NOTE:  GUI will be Slightly differ due to versions

# Iris Example Data Set



| | Features | | | | Class |
|---|---|---|---|---|---|
| | Sepal | | Petal | | Species |
| | *Length* | *Width* | *Length* | *Width* | |
| Pick flower 1 | | | | | 1 |
| Pick flower 2 | | | | | 2 |
| Pick flower 3 | | | | | 3 |
| Pick flower N | | | | | ???? |

**Instances**

# WEKA basics
## NOTE:  GUI will be Slightly differ due to versions

- API backs all functions of the CLI/GUI interfaces, can be easily used for your own project.

Program    Visualization    Tools    Help

Applications

WEKA
The University of Waikato

Explorer

Experimenter

KnowledgeFlow

Waikato Environment for Knowledge Analysis
Version 3.6.1
(c) 1999 – 2009
The University of Waikato
Hamilton, New Zealand

Simple CLI

# WEKA Explorer Tutorial Examples

- **Preprocess**
  - ➢ **Instance and Attribute Filters (Supervised and Unsupervised)**
- Classify
  - ➢ Bayes
- Cluster
  - ➢ Expectation Maximization
  - ➢ Hierarchical Clustering
- Associate
  - ➢ Apriori
- Select Attributes
  - ➢ Via clustering

# Preprocess

- File: CSV, ARFF*, ....
- Database: direct SQL access (useful)

# WEKA "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 "flat" files

@relation heart-disease-simplified

@attribute age numeric       **numeric attribute**

@attribute sex { female, male}       **nominal attribute**

@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 Explorer

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

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

## Filter

| Choose | **None** | | Apply |

## Current relation

Relation: None                    Attributes: None
Instances: None

## Selected attribute

Name: None                                        Type: None
Missing: None          Distinct: None             Unique: None

## Attributes

| Visualize All |

## Status

Welcome to the Weka Knowledge Explorer

| 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

| 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

21
16  16
13  14  14    15
10
7
5   6                    5       5
                    2       1
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: 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 |

## Attributes

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

Colour: class (Nom)          **Visualize All**

## Status
OK

**Log**   x 0
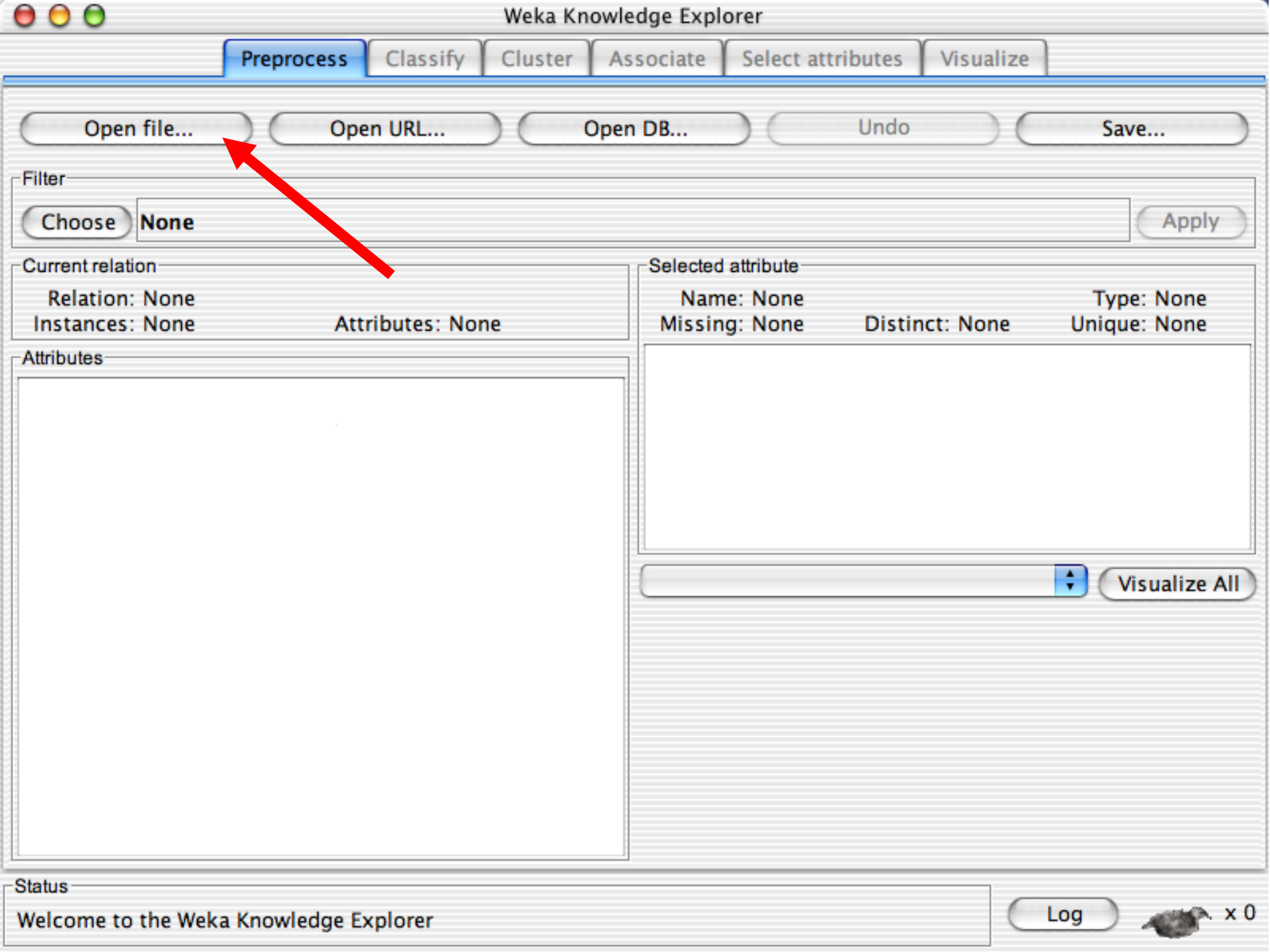
# Weka Knowledge Explorer

| 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

**Selected attribute**

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

**Attributes**

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

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

Colour: class (Nom)    Visualize All

50      50      50

**Status**

OK     Log    x 0

# Weka Knowledge Explorer

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

## Filter

Choose **None** | Apply

## Current relation

Relation: iris
Instances: 150    Attributes: 5

## 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 |

## Attributes

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

Colour: class (Nom)    Visualize All

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

## Status

OK

Log    x 0

# Weka Knowledge Explorer

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

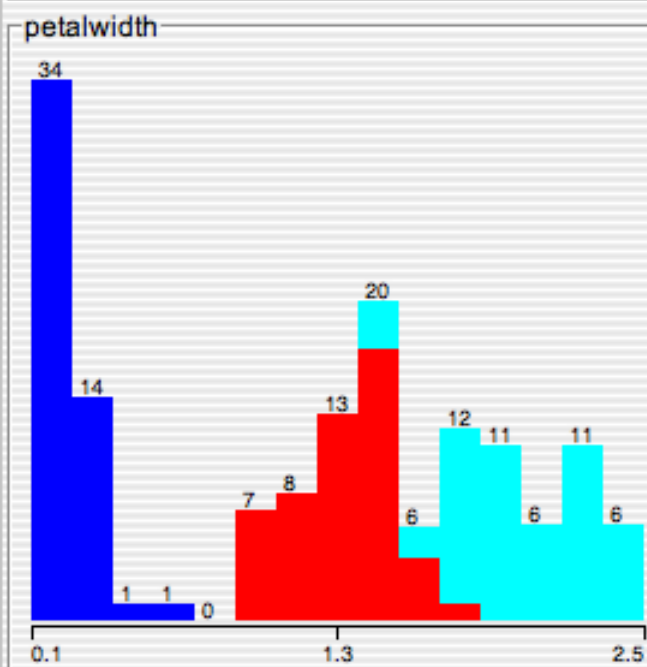**Open file...**   **Open URL...**   **Open DB...**   Undo   **Save...**

### Filter

**Choose** None   Apply

### Current relation
Relation: iris
Instances: 150   Attributes: 5

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

| Statistic | Value |
|-----------|-------|
| Minimum | 1 |
| Maximum | 6.9 |
| Mean | 3.759 |
| StdDev | 1.764 |

### Attributes

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

Colour: class (Nom)   **Visualize All**

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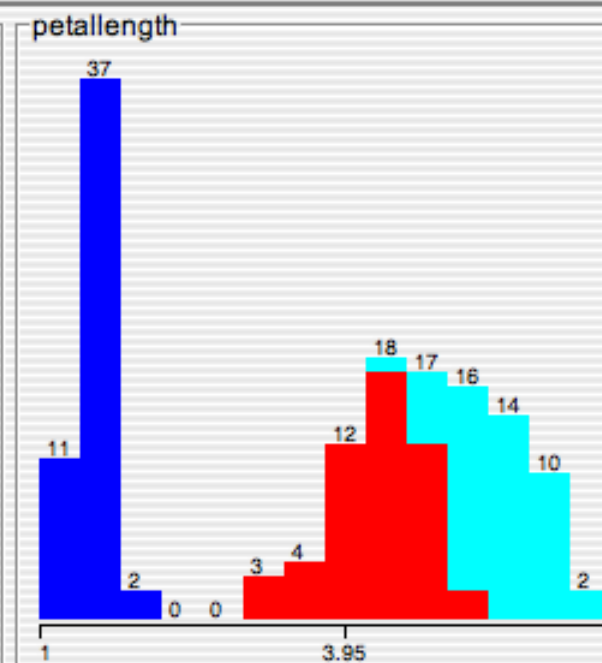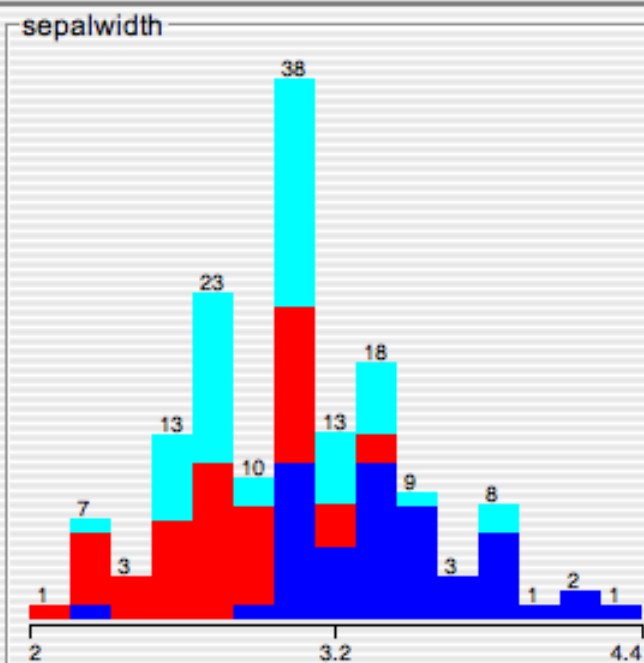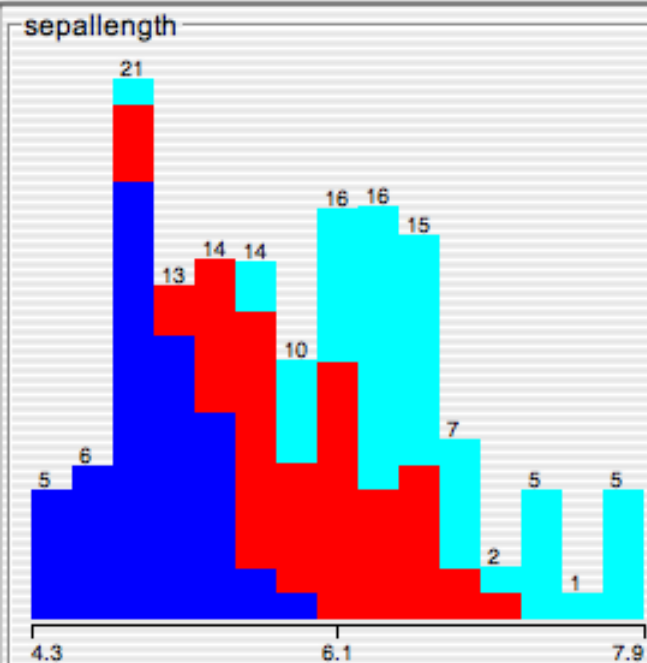1   3.95   6.9

### Status
OK   Log   x 0

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

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

**Filter**

Apply

- weka
  - filters
    - unsupervised
      - attribute
      - instance

**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

12   18   17   16   14

10

2   0   0   3   4                2   4

1          3.95          6.9

**Status**

OK

Log          x 0

# Weka Knowledge Explorer

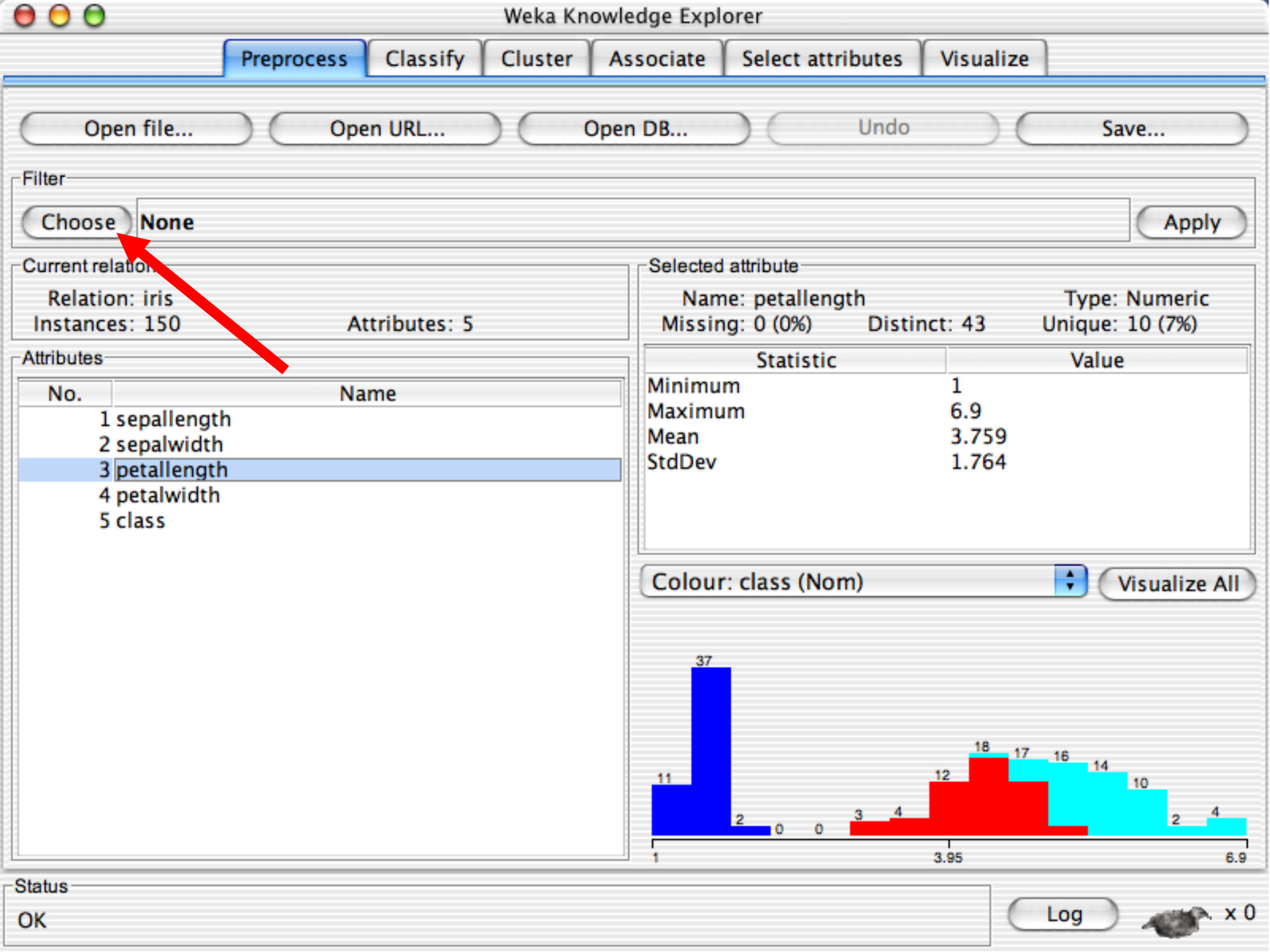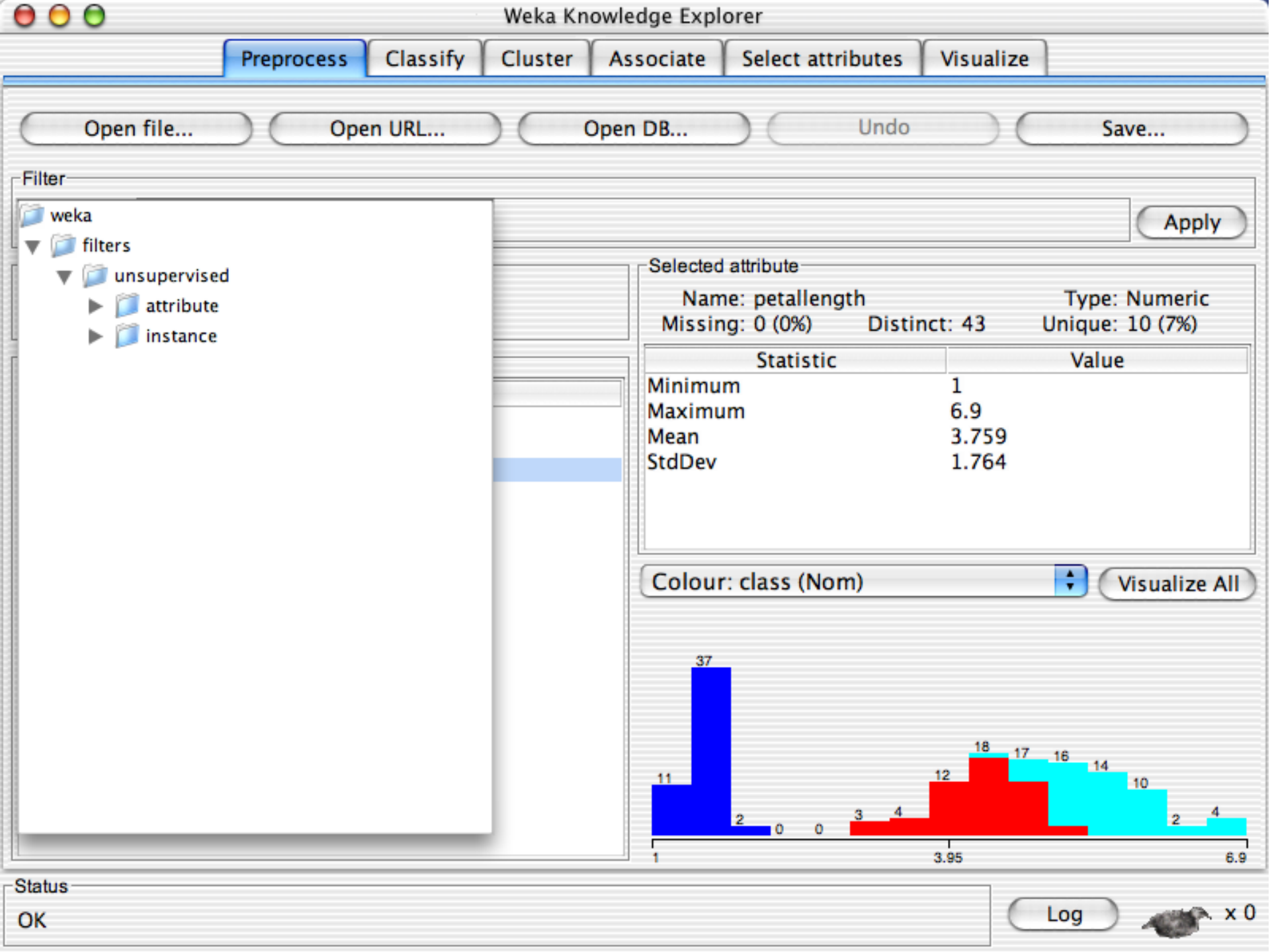| Preprocess | Classify | Cluster | Associate | Select attributes | Visualize |

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

## Filter

- weka
  - filters
    - unsupervised
      - attribute
      - instance

Apply

### 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
12  18  17  16  14
10
2   3   4           2   4
0   0

1         3.95         6.9

## Status

OK

Log    x 0

# Weka Knowledge Explorer

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

Apply

### 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

12   18   17   16   14   10

2   0   0   3   4   2   4

1          3.95          6.9
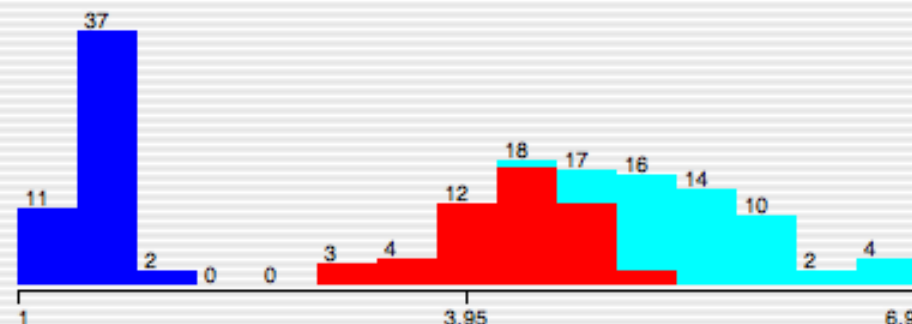
## Status

OK

Log          x 0

# Weka Knowledge Explorer

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

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

## Filter

Choose | **Discretize −B 10 −R first−last** | Apply

## Current relation

Relation: iris
Instances: 150    Attributes: 5

## 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 |

## Attributes

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

Colour: class (Nom) | Visualize All

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

1          3.95          6.9

## 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 | **Discretize** -B 10 -R first-last | Apply |
|---|---|---|

**Current relation**

Relation: iris
Instances: 150          Attributes: 5

**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 |

**Attributes**

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

Colour: class (Nom)          Visualize All

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

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

## Current relation

Relation: iris
Instances: 150          Attributes: | : Numeric

## Attributes

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

### weka.gui.GenericObjectEditor

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 |

Open... | Save... | OK | Cancel

: 10 (7%)

Visualize All

11

2
0          0          3   4                          10

2    4

1                              3.95                              6.9

## Status

OK                                              Log          x 0
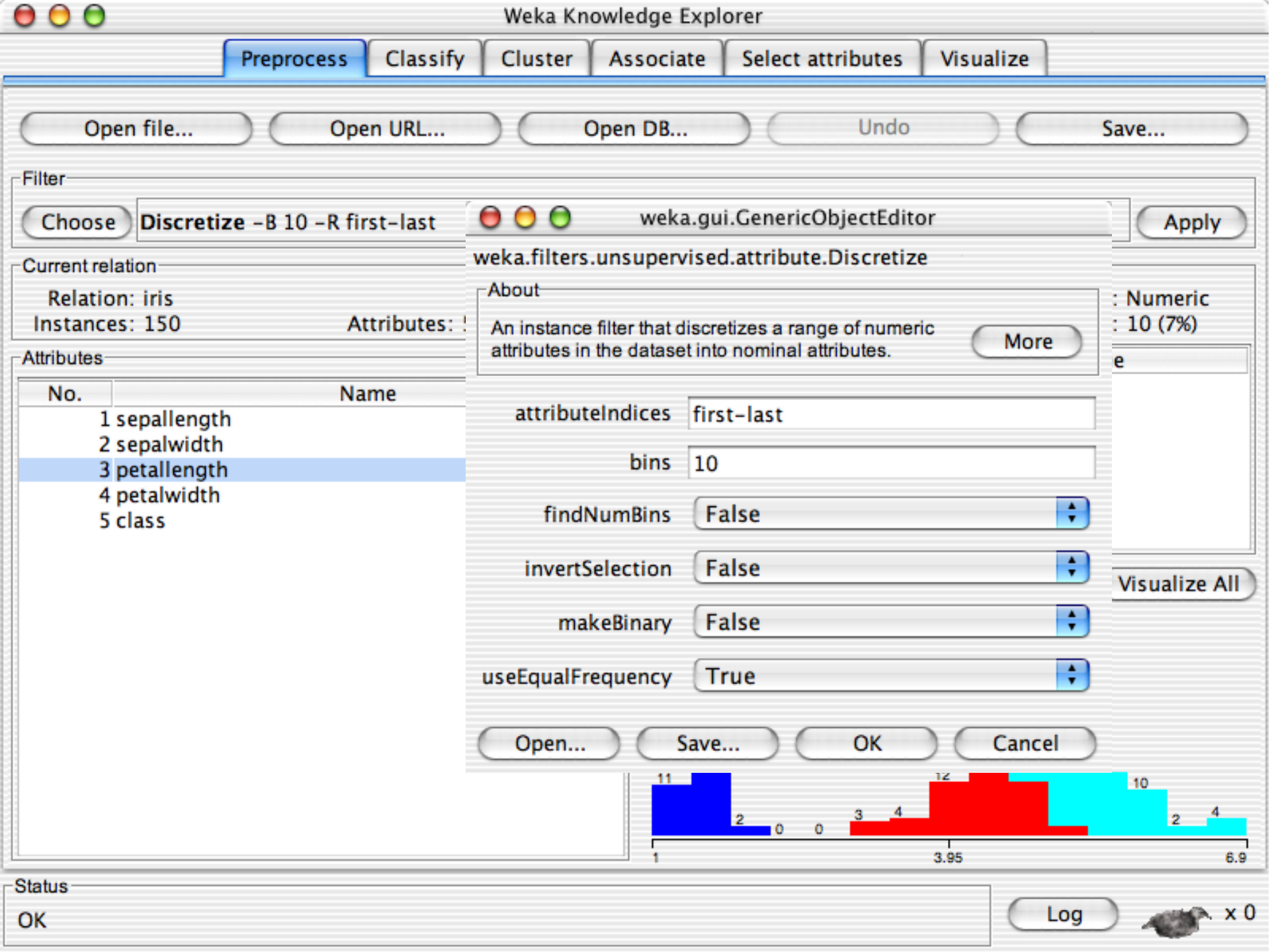
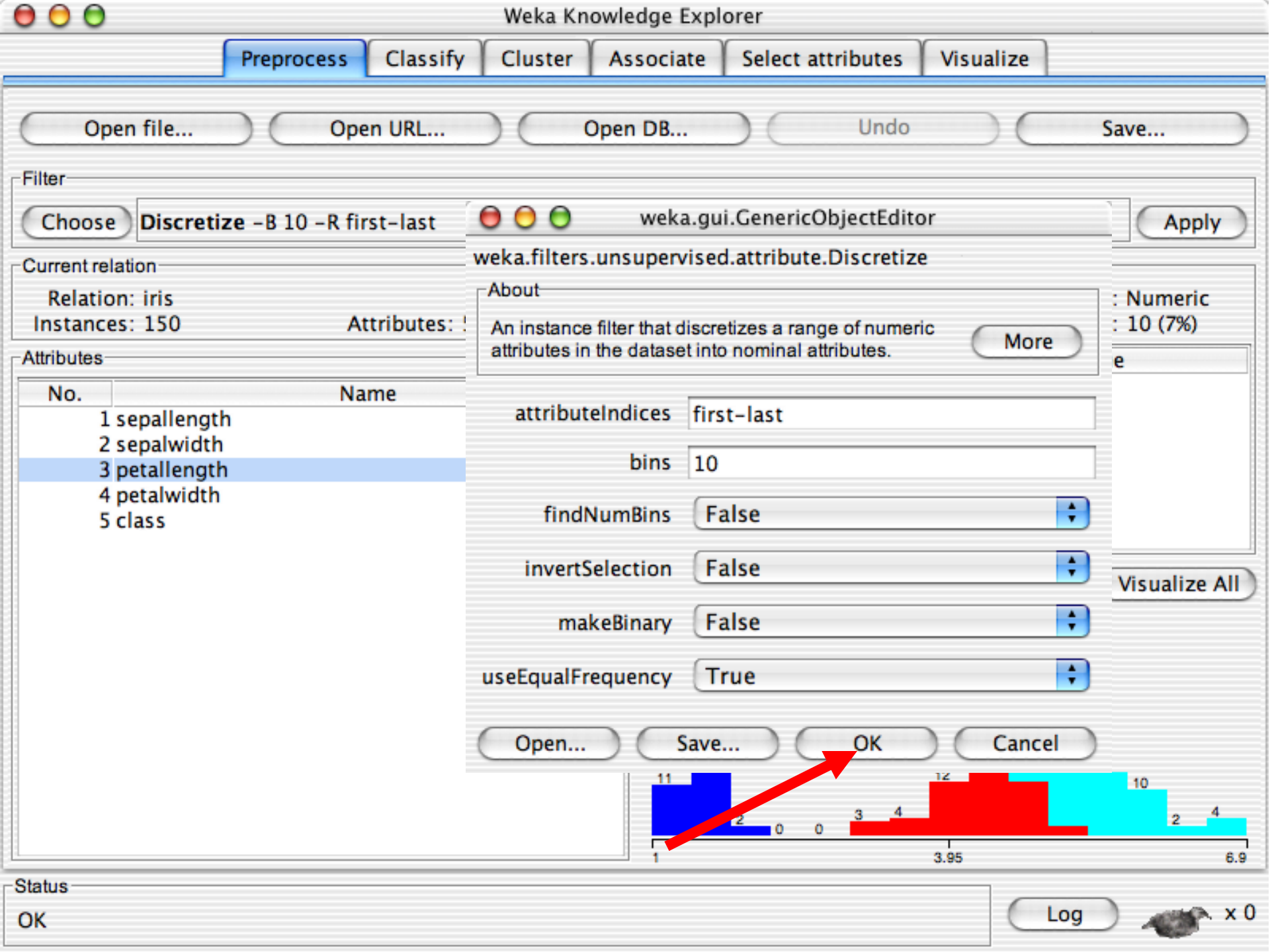# Weka Knowledge Explorer

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

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

## Filter

Choose    **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 |

### weka.gui.GenericObjectEditor

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 |

Open...    Save...    OK    Cancel

: Numeric
: 10 (7%)

Visualize All

11    12    10

2    3    4    2    4
0    0

1    3.95    6.9

## Status
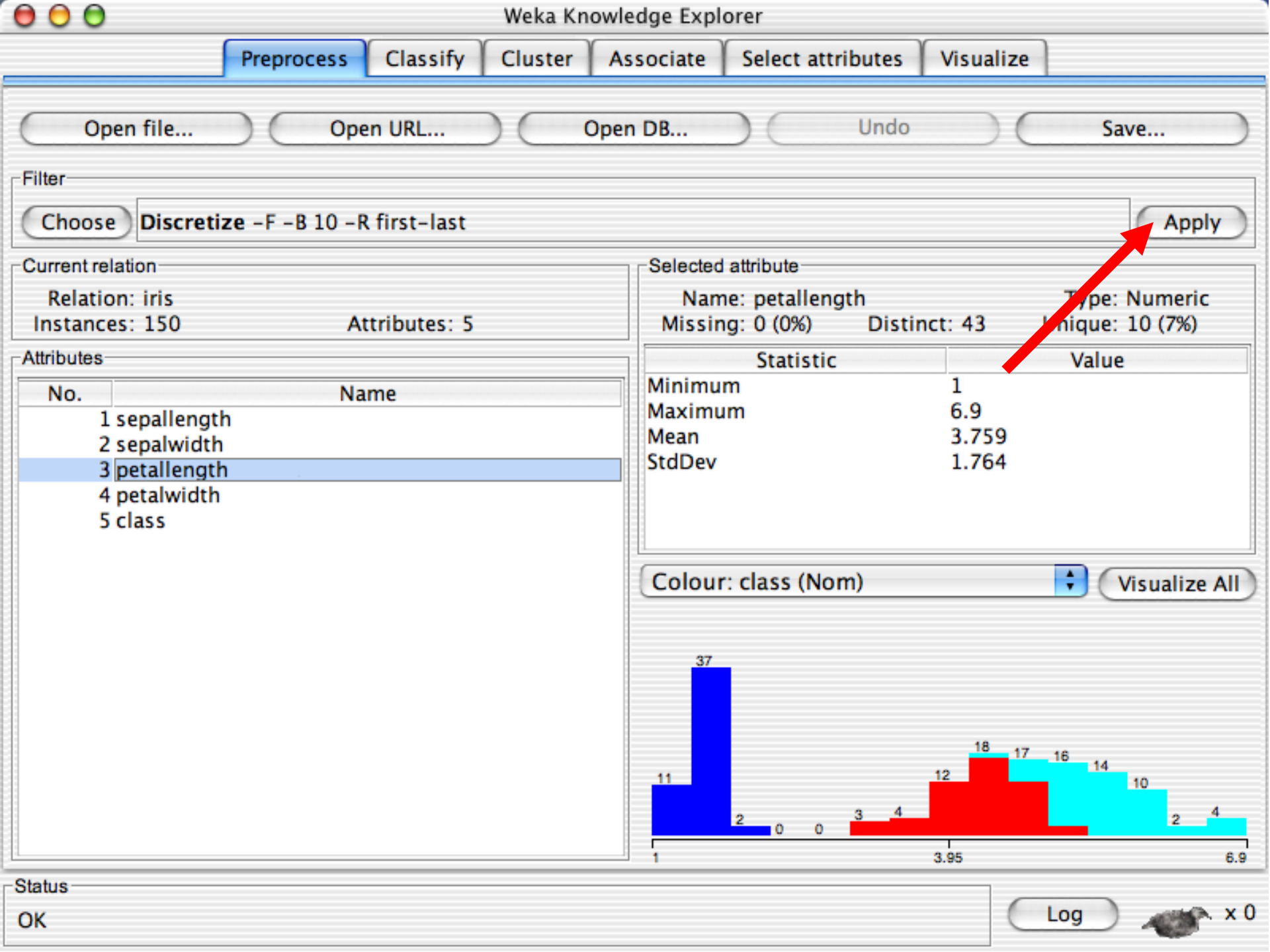
OK

Log    x 0

# Weka Knowledge Explorer

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

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

## Filter

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

## Current relation

Relation: iris
Instances: 150     Attributes:     : Numeric
    : 10 (7%)

## Attributes

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

### weka.gui.GenericObjectEditor

**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 | True |

Open... | Save... | OK | Cancel

Visualize All

11    12    10

2   0   0   3   4    2   4

1      3.95      6.9

## 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**  **Discretize -F -B 10 -R first-last**  Apply

## Current relation

Relation: iris
Instances: 150          Attributes: 5

## 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 |

## Attributes

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

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

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

## 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 |

## Attributes

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

Colour: class (Nom)    Visualize All

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

1    3.95    6.9

## 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 | **Discretize** -F -B 10 -R first-last | Apply

## Current relation

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

## 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 |

## Attributes

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

Colour: class (Nom) | Visualize All

## Status

OK

Log | x 0

# Preprocess: filters

- Choosing a filter
  - Supervised vs. Unsupervised
  - Attribute vs. Instance

- Supervised filters "require a class attribute"; unsupervised filters do not

- "meta-filters" can filter results from clustering and classification steps

# Preprocess: filter examples

|  | Instance | Attribute |
|---|---|---|
| **UNsupervised** | Resample | Discretize |
| **Supervised** | Resample | Discretize |

# Preprocess: <u>instance </u>filter example

- **Resampling**
  - <u>Unsupervised</u>
    - random % of the dataset

  - <u>Supervised</u>
    - takes the *class distribution* into account when generating a random sample
    - Can add bias towards a specific class value
    - Can specify maximum spread for rare/common class values

# Preprocess: <u>attribute</u> filter example

- **Discretize**
  - <u>Unsupervised</u>
    - K-Interval : simplest, can ensure small bin sizes

    - Proportional K-Interval : optimized for classification (Naïve Bayes)

  - <u>Supervised</u>
    - Entropy based
      - state of the art
      - computationally expensive

  - see *Chapter 7 of Data Mining by I. H. Witten and E. Frank*

# Preprocess: <u>attribute</u> filter favorites

- Finding and Discarding variables
  - RemoveUseless : cut using variation threshold

- Datatype Transforms
  - NumericToNominal & NominalToBinary
  - StringToWord : NLP

- Value transforms
  - Normalize
  - ReplaceMissingValues : with mean value from training data
  - AddExpression : any math expression (think R)

# WEKA Explorer Tutorial Examples

- Preprocess
  - Instance and Attribute Filters (Supervised and Unsupervised)
- **Classify**
  - ZeroR
  - Bayes
- Cluster
  - Expectation Maximization
  - Hierarchical Clustering
- Associate
  - Apriori
- Select Attributes
  - Via clustering

# 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

# Weka Knowledge Explorer

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

## Classifier

- weka
  - ▼ classifiers
    - ▶ bayes
    - ▶ functions
    - ▶ lazy
    - ▶ meta
    - ▶ misc
    - ▼ trees
      - ▶ adtree
      - DecisionStump
      - Id3
      - ▼ j48
        - J48
      - ▶ lmt
      - ▶ m5
      - RandomForest
      - RandomTree
      - REPTree
      - UserClassifier
    - ▶ rules

ifier output

## 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)

## Classifier output

## 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

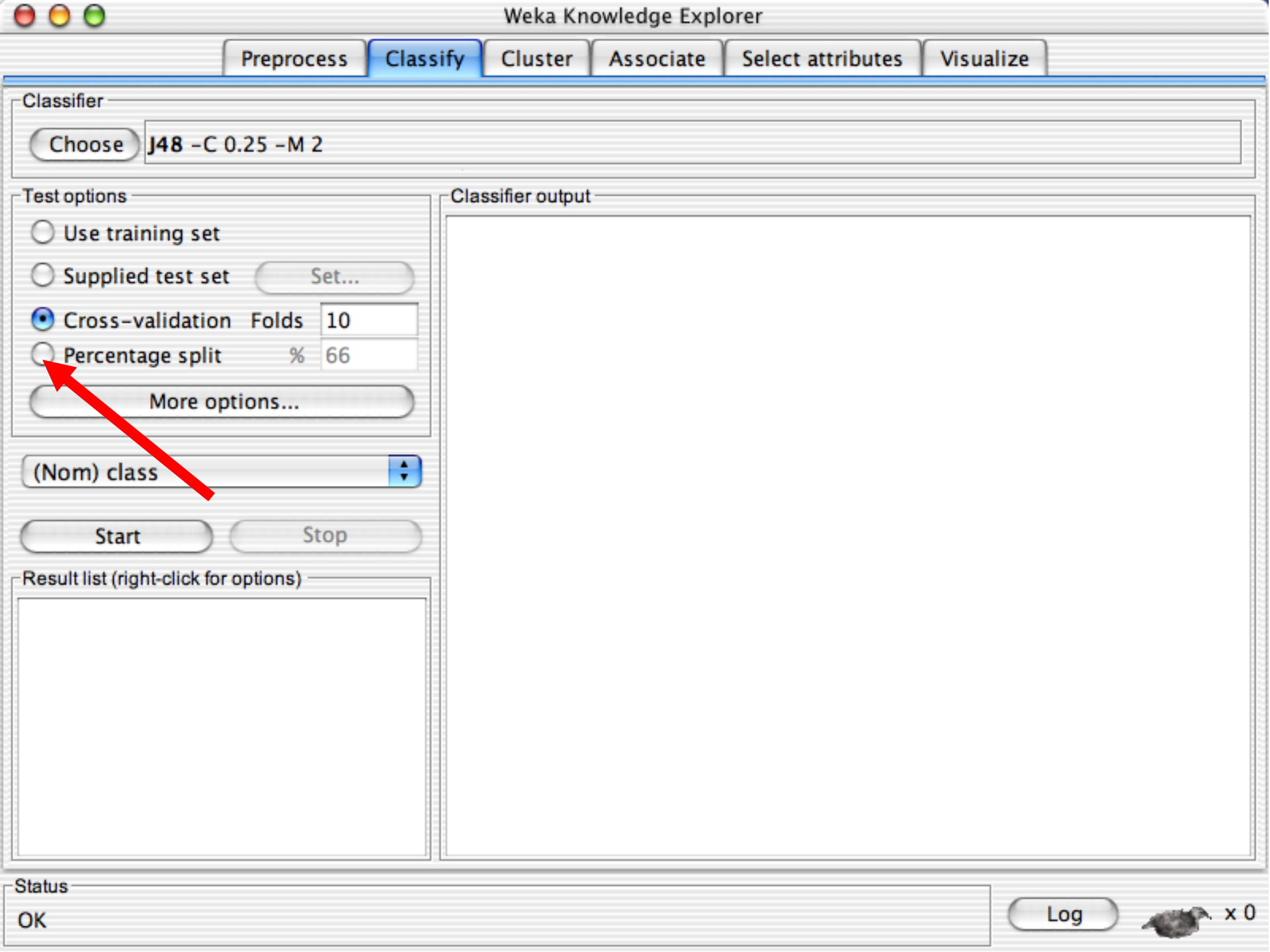**Start**    Stop

## Result list (right-click for options)

### weka.gui.GenericObjectEditor

weka.classifiers.trees.j48.J48

| | |
|---|---|
| binarySplits | False |
| confidenceFactor | 0.25 |
| minNumObj | 2 |
| numFolds | 3 |
| reducedErrorPruning | False |
| saveInstanceData | False |
| subtreeRaising | True |
| unpruned | False |
| useLaplace | False |

Open...   Save...   OK   Cancel

## 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
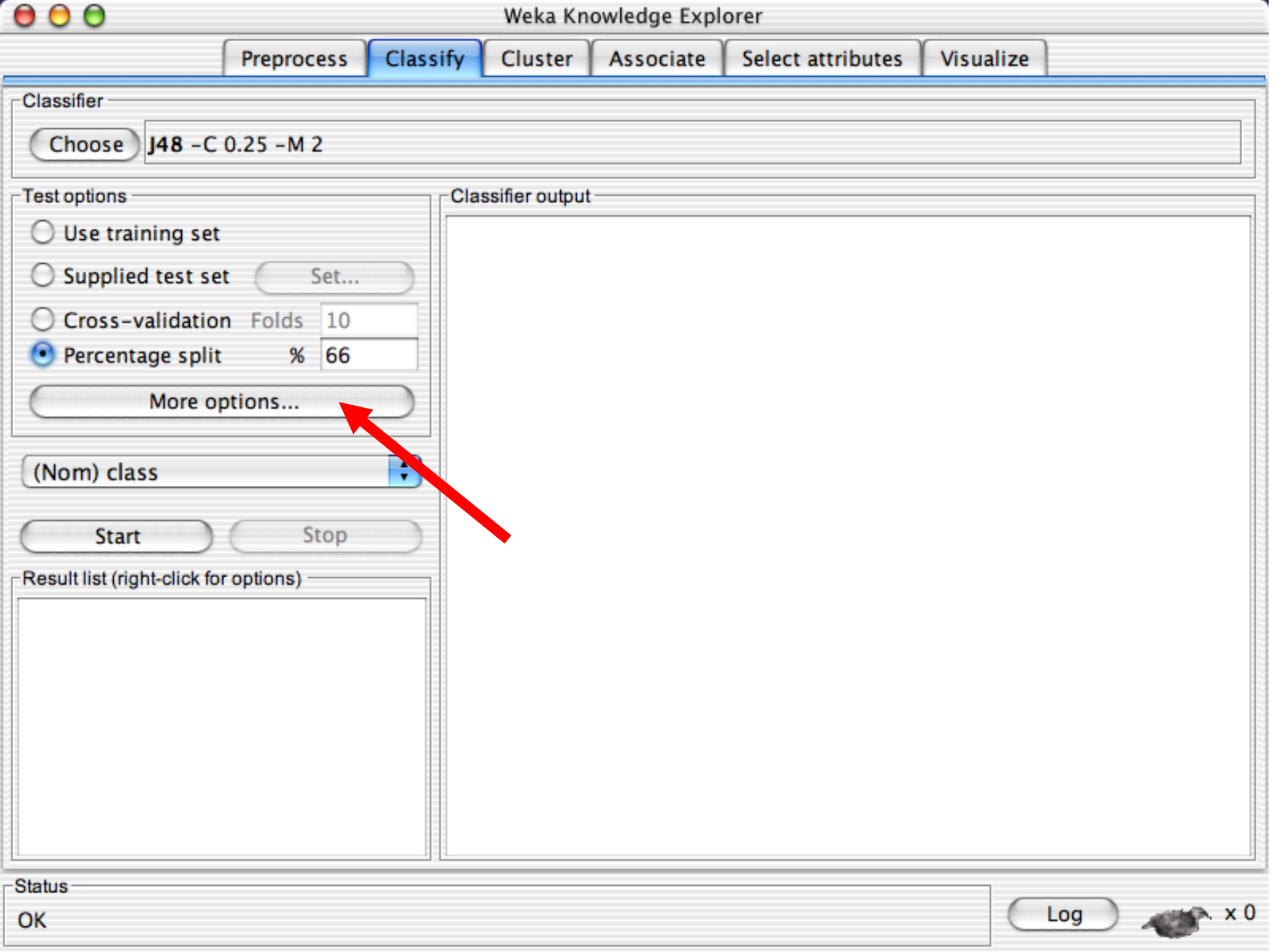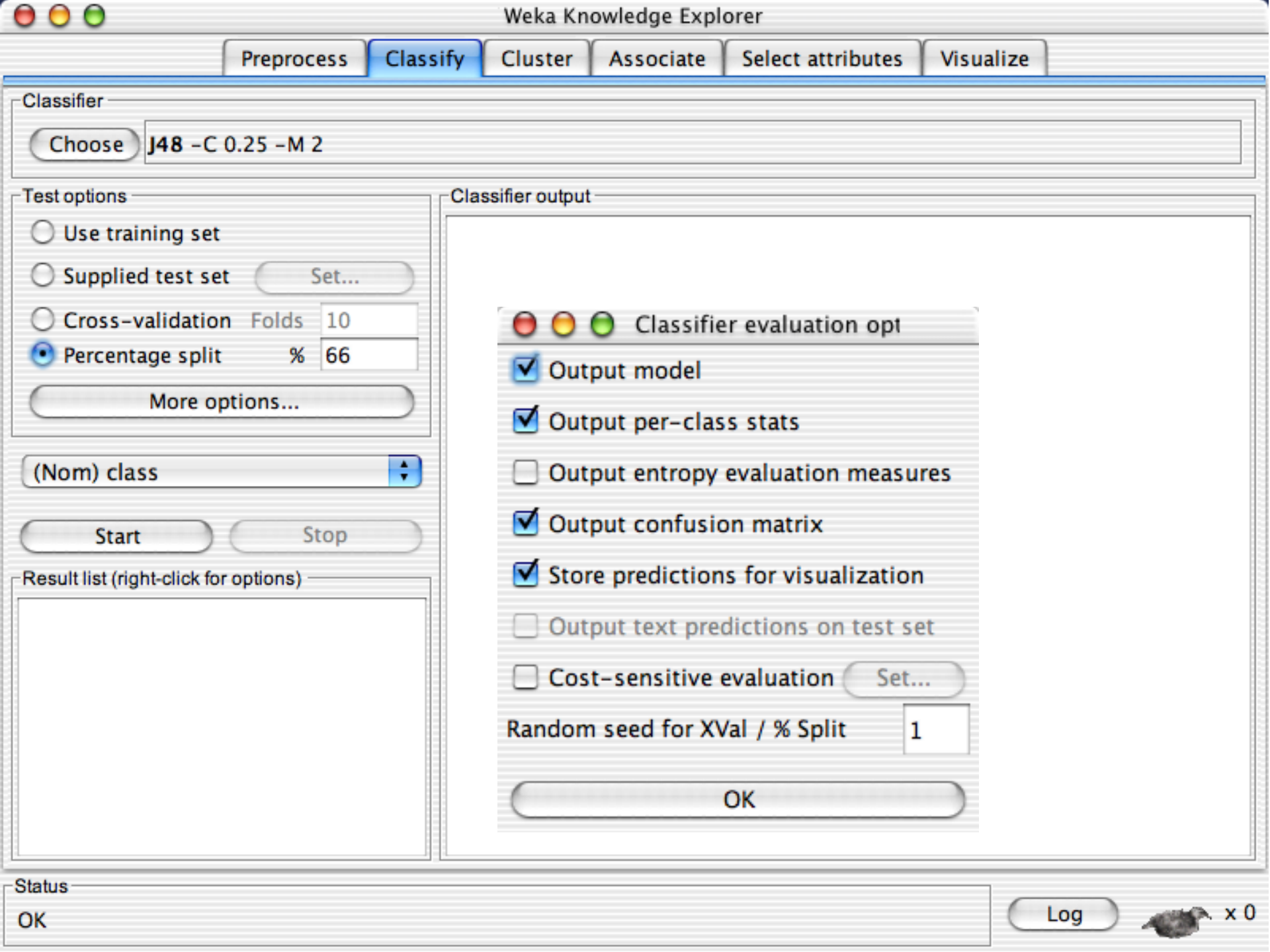
Start | Stop

Result list (right-click for options)

---

### weka.gui.GenericObjectEditor

weka.classifiers.trees.j48.J48

| | |
|---|---|
| binarySplits | False |
| confidenceFactor | 0.25 |
| minNumObj | 2 |
| numFolds | 3 |
| reducedErrorPruning | False |
| saveInstanceData | False |
| subtreeRaising | True |
| unpruned | False |
| useLaplace | False |

Open... | Save... | OK | Cancel

---

## 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)

## Classifier output

## 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)**

**Classifier output**

**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)

## Classifier output

## 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

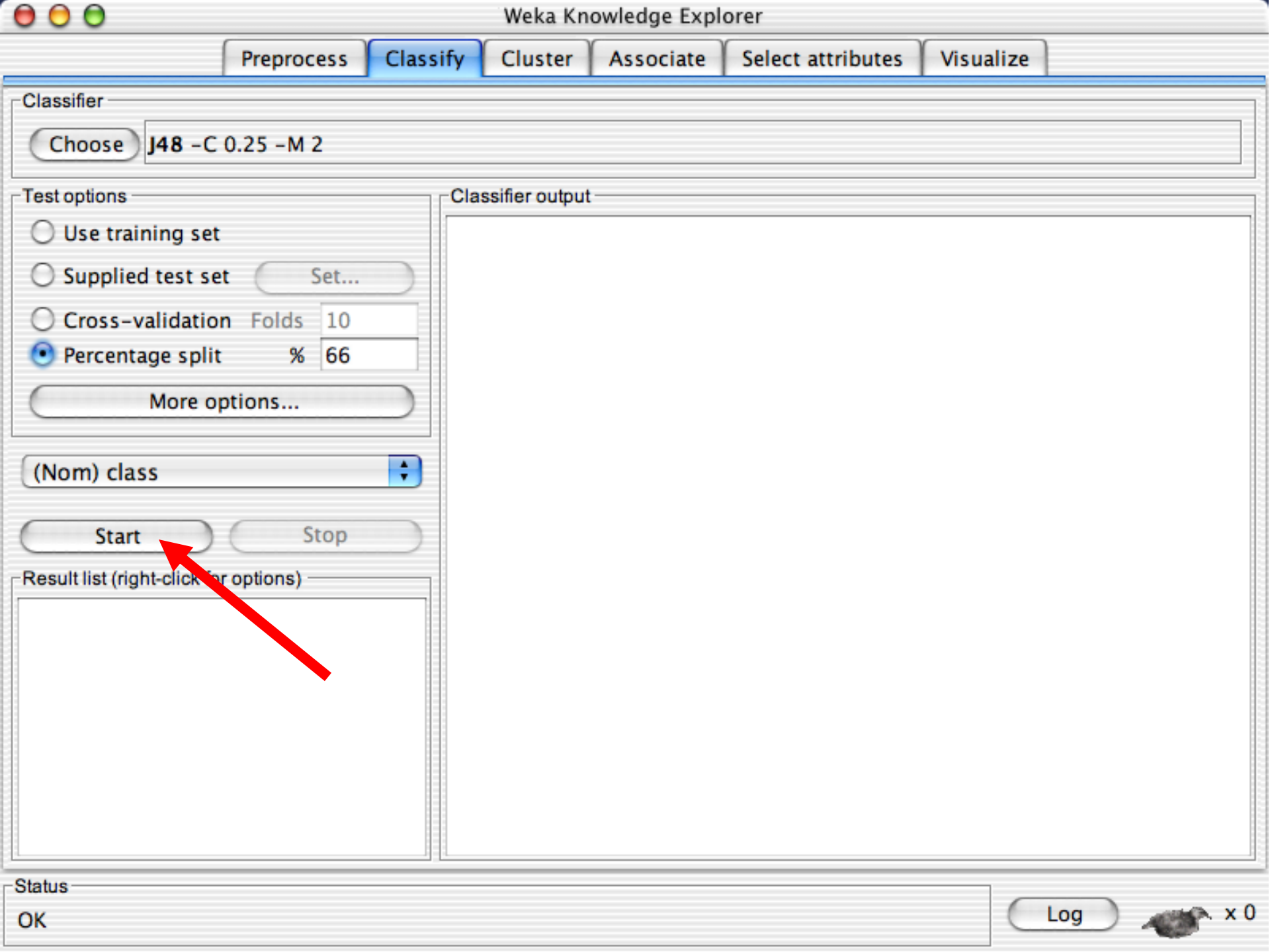## Result list (right-click for options)

## Classifier output

### Classifier evaluation options

☑ Output model

☑ Output per-class stats

☐ Output entropy evaluation measures

☑ Output confusion matrix

☑ Store predictions for visualization

☐ Output text predictions on test set

☐ Cost-sensitive evaluation    Set...

Random seed for XVal / % Split    1

OK

## 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)

## Classifier output

## 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 ]

## Classifier output

## Result list (right-click for options)

## 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

```
=== 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
```
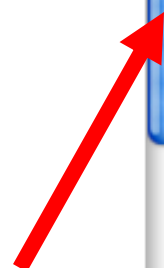
## Status

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

```
=== 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
```

## 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

**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
```
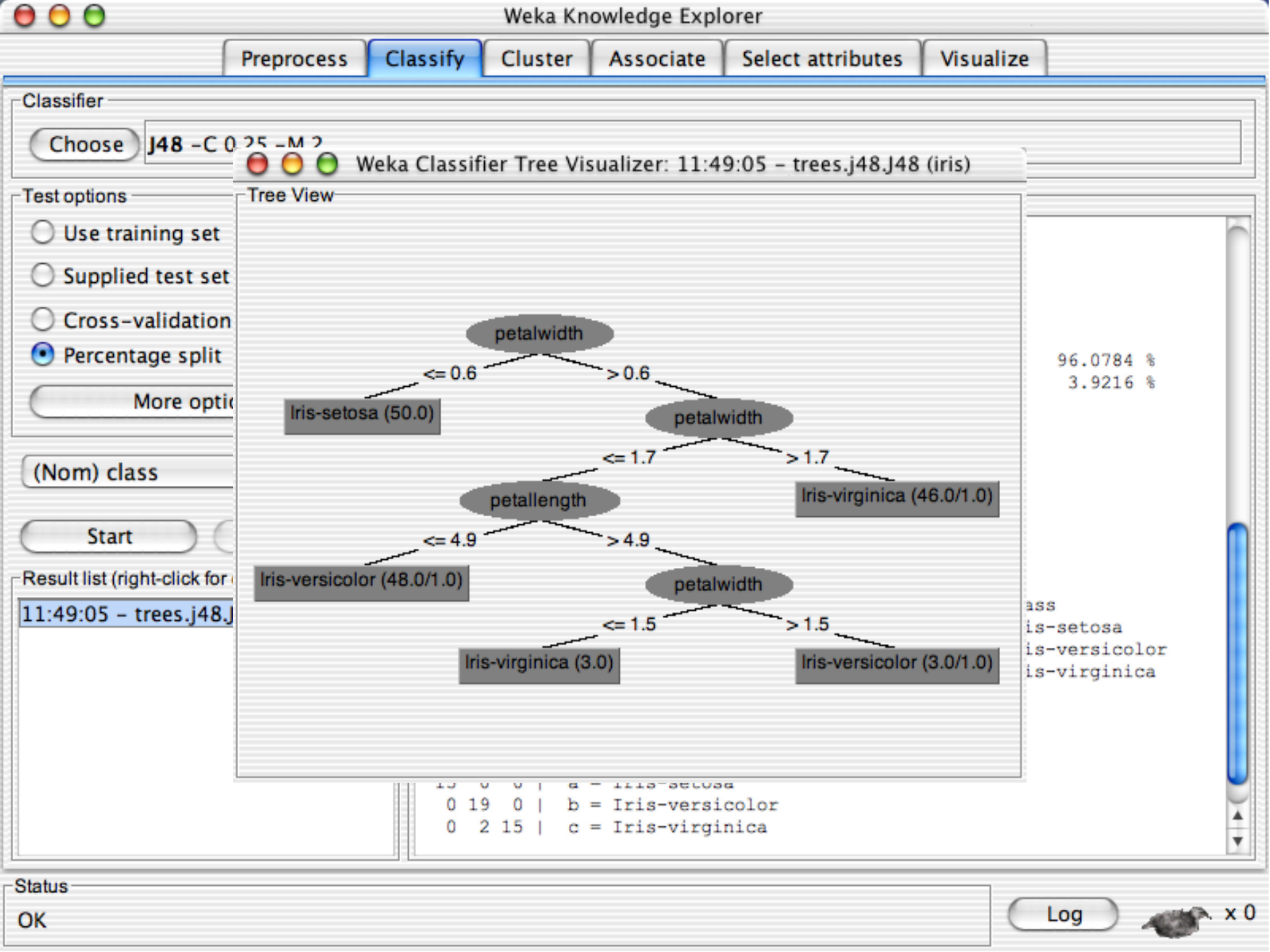
## Status

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 ===
```

|  | 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 classifer errors
**Visualize tree**
Visualize margin curve
Visualize threshold curve    ▶
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)

### Tree View

## Test options

○ Use training set

○ Supplied test set

○ Cross-validation

● Percentage split

More optic

(Nom) class

Start

Result list (right-click for

11:49:05 – trees.j48.J



96.0784 %
3.9216 %

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

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 classifer errors**
Visualize tree
Visualize margin curve
Visualize threshold curve          ▶
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

Result list (right-click for ...

11:49:05 – trees.j48.J...

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

X: petallength (Num)    Y: petalwidth (Num)

Colour: class (Nom)    Select Instance

Reset | Clear | Save    Jitter

96.0784 %
3.9216 %

### Plot: iris_predicted



### Class colour

Iris-setosa  Iris-versicolor  Iris-virginica

    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

# Classify : Bayes

Various Bayesian network classifier learning algorithms are implemented in Weka [12]. This note provides some user documentation and implementation details.

Summary of main capabilities:

- Structure learning of Bayesian networks using various hill climbing (K2, B, etc) and general purpose (simulated annealing, tabu search) algorithms.

- Local score metrics implemented; Bayes, BDe, MDL, entropy, AIC.

- Global score metrics implemented; leave one out cv, k-fold cv and cumulative cv.

- Conditional independence based causal recovery algorithm available.

- Parameter estimation using direct estimates and Bayesian model averaging.

- GUI for easy inspection of Bayesian networks.

- Part of Weka allowing systematic experiments to compare Bayes net performance with general purpose classifiers like C4.5, nearest neighbor, support vector etc.

# Classify: Learning a Bayes Network Structure



CAI3024N-ADA-JJT-WK7

# Classify using Naïve Bayes

- Naïve Bayes                   (standard )
- Multinominal Naïve Bayes     (text classification)

- Hidden Naïve Bayes, others, ….

# Classify using Naïve Bayes

1. Pick Discretization method* (hardest part)
2. Pick a class to predict
3. Run the classifier

# Classify Text using Naïve Bayes Multinominal

- Framingham dataset contains text descriptions for each of the 21k+ variables

- I wrote a parallelized NLP program to calculate inverse word frequencies and score variable pairs ( 2 days )

- Highest scoring pairs were suggested for merger to reduce the variable space (curse of dimensionality)

# Classify Text using Naïve Bayes Multinominal

*What I should have done....*

① Train an NBC to learn from small set of labeled cases

② Apply NBC to unlabeled data using Expectation Maximization with class probabilities (expectation step)

③ Retrain NBC using the labels for all the data

④ Repeat until convergence

# WEKA Explorer Tutorial Examples

- Preprocess
  - Instance and Attribute Filters (Supervised and Unsupervised)
- Classify
  - Bayes
- **Cluster**
  - **Expectation Maximization**
  - **Hierarchical Clustering**
- Associate
  - Apriori
- Select Attributes
  - Via clustering

# Cluster Algorithm Examples

- Expectation Maximization (EM)

- Hierarchical Clustering (cobweb)

- *Note: Weka provides many more clustering methods*

# Cluster with Expectation Maximization

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

**Clusterer**

[ Choose ]  **EM** -I 100 -N -1 -M 1.0E-6 -S 100

**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)**

04:16:46 - EM

**Clusterer output**

```
                       Cluster
Attribute                 0       1       2       3
                       (0.32) (0.33)   (0.2)  (0.14)
======================================================
sepallength
  mean                  5.897   5.006  6.9426  6.1304
  std. dev.            0.5279  0.3489   0.498  0.2943

sepalwidth
  mean                 2.7519   3.418  3.1103  2.8088
  std. dev.            0.3103  0.3772  0.2952  0.2361

petallength
  mean                 4.2267   1.464  5.8559  5.0993
  std. dev.             0.445  0.1718  0.4626  0.2462

petalwidth
  mean                 1.3134   0.244  2.1495  1.8254
  std. dev.            0.1864  0.1061   0.232  0.2152

class
  Iris-setosa               1      51       1       1
  Iris-versicolor     48.1125       1  1.0182  3.8693
  Iris-virginica       2.0983       1 31.0375 19.8641
  [total]             51.2108      53 33.0557 24.7335
Clustered Instances

0      48 ( 32%)
1      50 ( 33%)
2      29 ( 19%)
3      23 ( 15%)
```

# Cluster with cobweb (hierarchical clustering)

# WEKA Explorer Tutorial Examples

- Preprocess
  - Instance and Attribute Filters (Supervised and Unsupervised)
- Classify
  - Bayes
- Cluster
  - Expectation Maximization
  - Hierarchical Clustering
- **Associate**
  - **Apriori**
- Select Attributes
  - Via clustering

# Associate

- Quick scan for association rules
  - see "Fast Algorithms for Mining Association Rules in Large Databases"



| 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 –c –1

Start | Stop

Result list (right-click fo

04:47:46 – HotSpot
05:11:10 – Apriori

**Associator output**

```
Minimum support: 0.1 (15 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 18

Generated sets of large itemsets:

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

Size of set of large itemsets L(2): 15

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

Best rules found:

 1. petalwidth='(-inf-0.34]' 41 ==> class=Iris-setosa 41    conf:(1)
 2. petallength='(-inf-1.59]' 37 ==> class=Iris-setosa 37    conf:(1)
 3. petallength='(-inf-1.59]' petalwidth='(-inf-0.34]' 33 ==> class=Iris-setosa 33    conf:(1)
 4. petalwidth='(1.06-1.3]' 21 ==> class=Iris-versicolor 21    conf:(1)
 5. petallength='(5.13-5.72]' 18 ==> class=Iris-virginica 18    conf:(1)
 6. sepallength='(4.66-5.02]' petalwidth='(-inf-0.34]' 17 ==> class=Iris-setosa 17    conf:(1)
 7. sepalwidth='(2.96-3.2]' class=Iris-setosa 16 ==> petalwidth='(-inf-0.34]' 16    conf:(1)
 8. sepalwidth='(2.96-3.2]' petalwidth='(-inf-0.34]' 16 ==> class=Iris-setosa 16    conf:(1)
 9. petallength='(3.95-4.54]' 26 ==> class=Iris-versicolor 25    conf:(0.96)
10. petalwidth='(1.78-2.02]' 23 ==> class=Iris-virginica 22    conf:(0.96)
```

# WEKA Explorer Tutorial Examples

- Preprocess
  - Instance and Attribute Filters (Supervised and Unsupervised)
- Classify
  - Bayes
- Cluster
  - Expectation Maximization
  - Hierarchical Clustering
- Associate
  - Apriori
- **Select Attributes**
  - **Via clustering**

# Select Attributes Using a Classifier

# Select Attributes using PCA

Attribute Evaluator

[Choose] **PrincipalComponents** -R 0.95 -A 5

Search Method

[Choose] **Ranker** -T -1.7976931348623157E308 -N -1

Attribute Selection Mode

◉ Use full training set

○ Cross-validation    Folds  10

Seed  1

(Nom) class  ⬍

[Start]  [Stop]

Result list (right-click for options)

04:52:12 - GreedyStepwise + ClassifierSubse
04:58:49 - Ranker + PrincipalComponents

Attribute selection output

```
=== Run information ===

Evaluator:    weka.attributeSelection.PrincipalComponents -R 0.95 -A 5
Search:       weka.attributeSelection.Ranker -T -1.7976931348623157E308 -N -1
Relation:     iris
Instances:    150
Attributes:   5
              sepallength
              sepalwidth
              petallength
              petalwidth
              class
Evaluation mode:    evaluate on all training data


=== Attribute Selection on all input data ===

Search Method:
        Attribute ranking.

Attribute Evaluator (unsupervised):
        Principal Components Attribute Transformer

Correlation matrix
  1     -0.11   0.87   0.82
 -0.11   1     -0.42  -0.36
  0.87  -0.42   1      0.96
  0.82  -0.36   0.96   1


eigenvalue        proportion      cumulative
 2.91082           0.7277          0.7277       0.581petallength+0.566petalwidth+0.522sepallength-0.263sepalwidth
 0.92122           0.23031         0.95801      -0.926sepalwidth-0.372sepallength-0.065petalwidth-0.021petallength

Eigenvectors
 V1      V2
 0.5224 -0.3723 sepallength
-0.2634 -0.9256 sepalwidth
 0.5813 -0.0211 petallength
 0.5656 -0.0654 petalwidth

Ranked attributes:
 0.2723  1 0.581petallength+0.566petalwidth+0.522sepallength-0.263sepalwidth
 0.042   2 -0.926sepalwidth-0.372sepallength-0.065petalwidth-0.021petallength

Selected attributes: 1,2 : 2
```

# **WEKA Tutorial Summary**

- **Preprocess**
  - ➤ Prepare datasets instances and attributes before analysis

- **Classify**
  - ➤ Pick a instance and predict the class
    - ✦ Iris : Pick a flower and use the attributes to predict species
    - ✦ Medicine: pick a patient and use the genes to predict cancer status

- **Cluster**
  - ➤ Group instances together (flowers, breast cancer cases, etc)

- **Associate**
  - ➤ Discover relationships between variables in your dataset

# References

- Data Mining: Practical Machine Learning Tools and Techniques

- Data Mining (I. H. Witten and E. Frank)

- WEKA Exploratory Tool for Data Mining

- Bayesian Network Classifiers in Weka (Bouckaert)

- COC131 Data Mining – Clustering (Sykora)

- Fast, Correct Multithreaded Programs in Java (Gilbert)

- R. Agrawal, R. Srikant: Fast Algorithms for Mining Association Rules in Large Databases. In: 20th International Conference on Very Large Data Bases, 478-499, 1994.

- WEKA Wiki
  http://weka.wikispaces.com/

- Graphical User Interface
  http://prdownloads.sourceforge.net/weka/weka.ppt

# References

- **OL- Ebook**