

SIMBig14



INTRODUCTION TO WEKA

Université de Montpellier 2

Hugo Alatrasta-Salas : hugo.alatrasta-salas@teledetection.fr

Introduction

WEKA

- Gallirallus australis : Endemic bird (New Zealand)



Characteristics

- Waikato university
- Weka is a collection of machine learning algorithms for data mining tasks
- Weka contains tools for data pre-processing, classification, regression, clustering, association rules, and visualization.
- Under GPL license

Links

- ◎ <http://www.cs.waikato.ac.nz/ml/weka/>

- ◎ [http://transact.dl.sourceforge.net/
sourceforge/weka/
WekaManual-3.6.0.pdf](http://transact.dl.sourceforge.net/sourceforge/weka/WekaManual-3.6.0.pdf)

How to run Weka?

- Using the icon

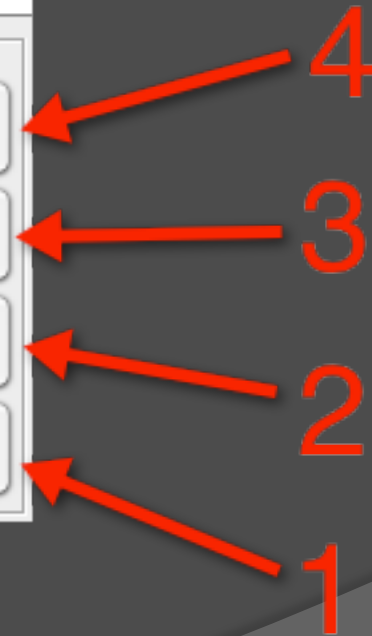
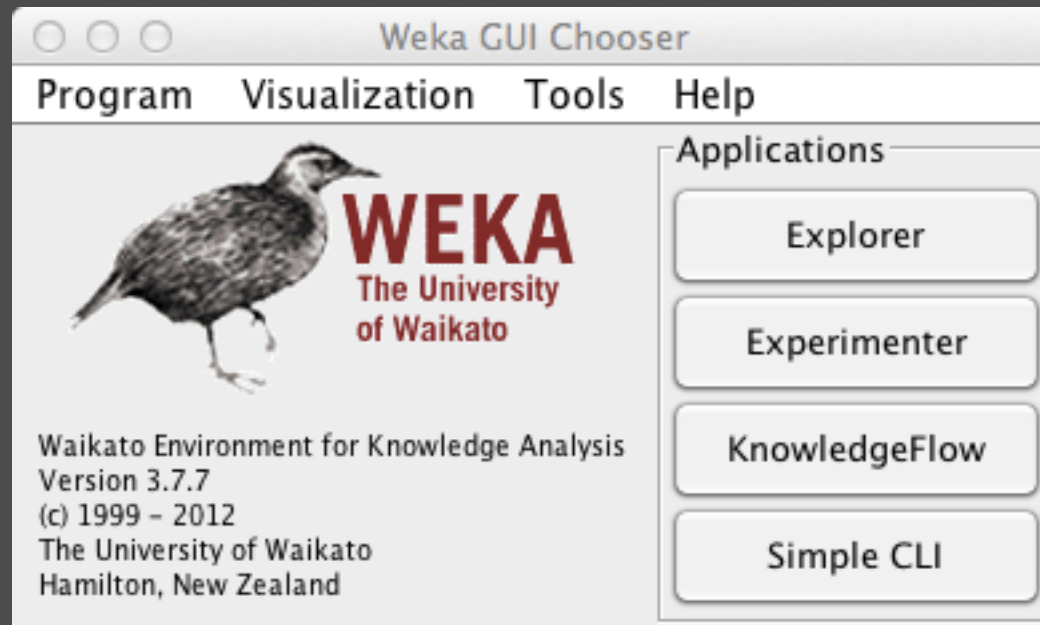


- Using the command line

```
java -Xmx1024m -jar weka.jar
```

Interface

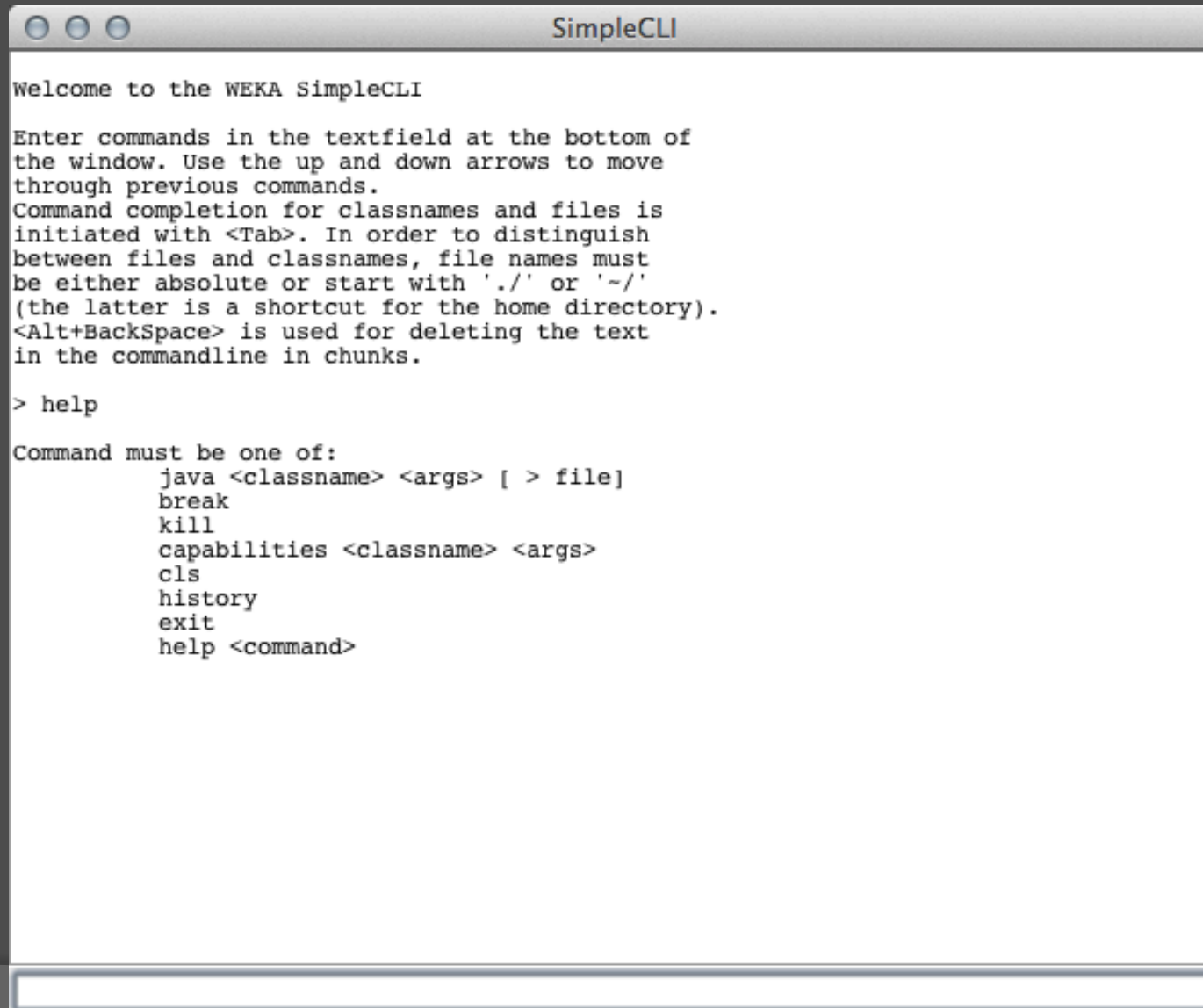
Main interface



Simple CLI (1)

- ⦿ Support all operation proposed by WEKA
- ⦿ E.g.
 - java <class><param>
 - break
 - kill
 - cls
 - exit
 - help <command>
 - ...

Simple CLI (2)



```
SimpleCLI

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.
Command completion for classnames and files is
initiated with <Tab>. In order to distinguish
between files and classnames, file names must
be either absolute or start with './' or '~/
(the latter is a shortcut for the home directory).
<Alt+BackSpace> is used for deleting the text
in the commandline in chunks.

> help

Command must be one of:
    java <classname> <args> [ > file]
    break
    kill
    capabilities <classname> <args>
    cls
    history
    exit
    help <command>
```

Knowledge Flow (1)

- ⦿ Alternative to the Explorer as a graphical front end
- ⦿ Intuition:

The user can select WEKA components from a tool bar, place them on a layout canvas and connect them together in order to form a knowledge flow for processing and analyzing data.

Knowledge Flow (2)

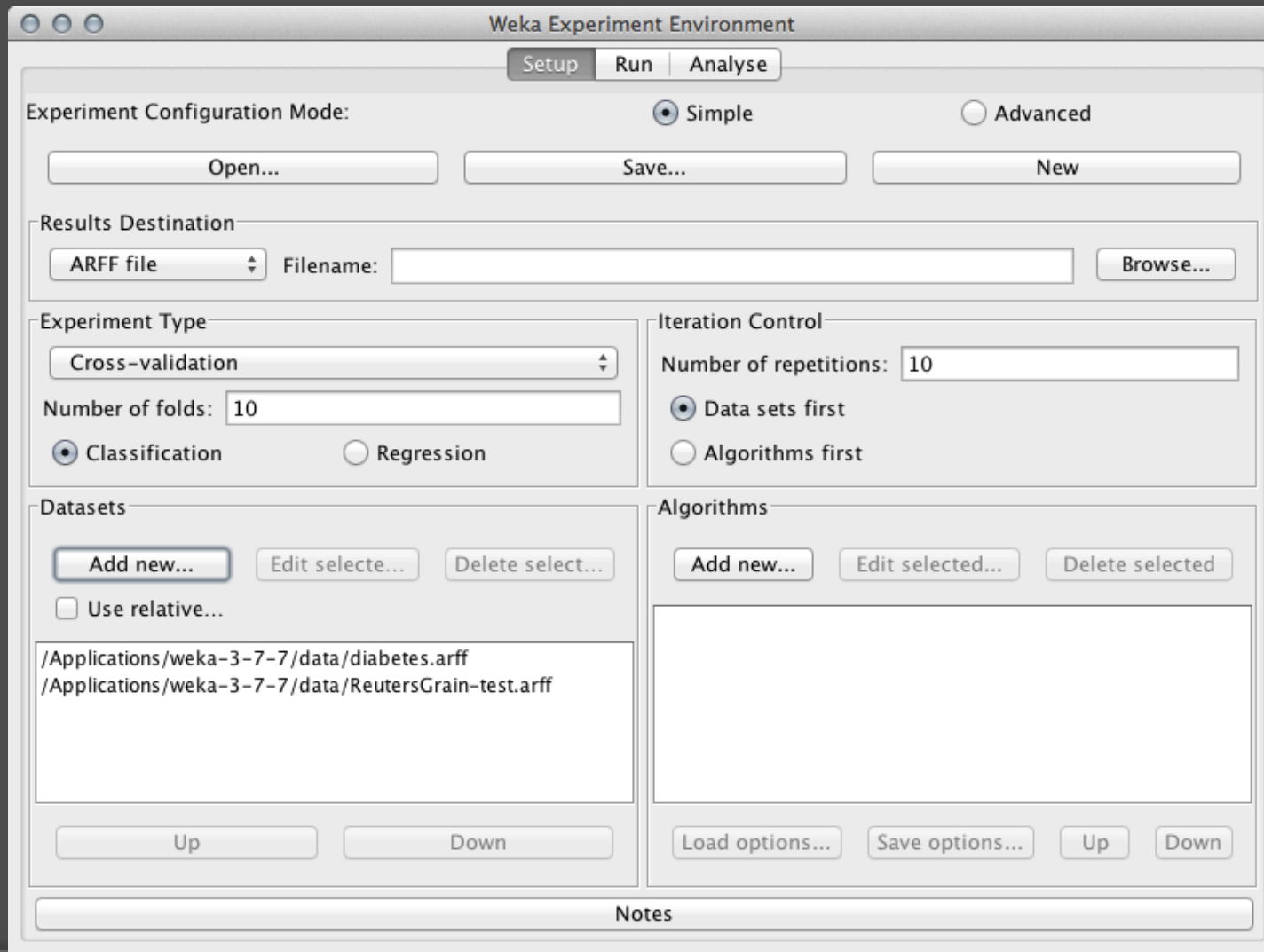
The screenshot displays the Weka KnowledgeFlow Environment interface. The title bar reads "Weka KnowledgeFlow Environment". Below the title bar, there is a toolbar with various icons for file operations and execution. The main workspace, titled "Untitled1", contains a workflow diagram. On the left, a "Design" panel lists several categories: DataSources, DataSinks, Filters, Classifiers, Clusterers, Associations, Evaluation, Visualization, and Tools. The workflow in the center consists of two nodes: "ArffLoader" on the left and "Discretize" on the right. A red arrow labeled "dataSet" connects the two nodes. At the bottom of the interface, there is a "Status" and "Log" section. The status table shows the following information:

Component	Parameters	Time	Status
[KnowledgeFlow]		0:5:35	OK.

Experimenter (1)

- ⦿ Experimenter makes it easy to compare the performance of different learning schemes
- ⦿ For classification and regression problems
- ⦿ Results can be written into file or database
- ⦿ Evaluation options: cross-validation, learning curve, hold-out
- ⦿ Can also iterate over different parameter settings

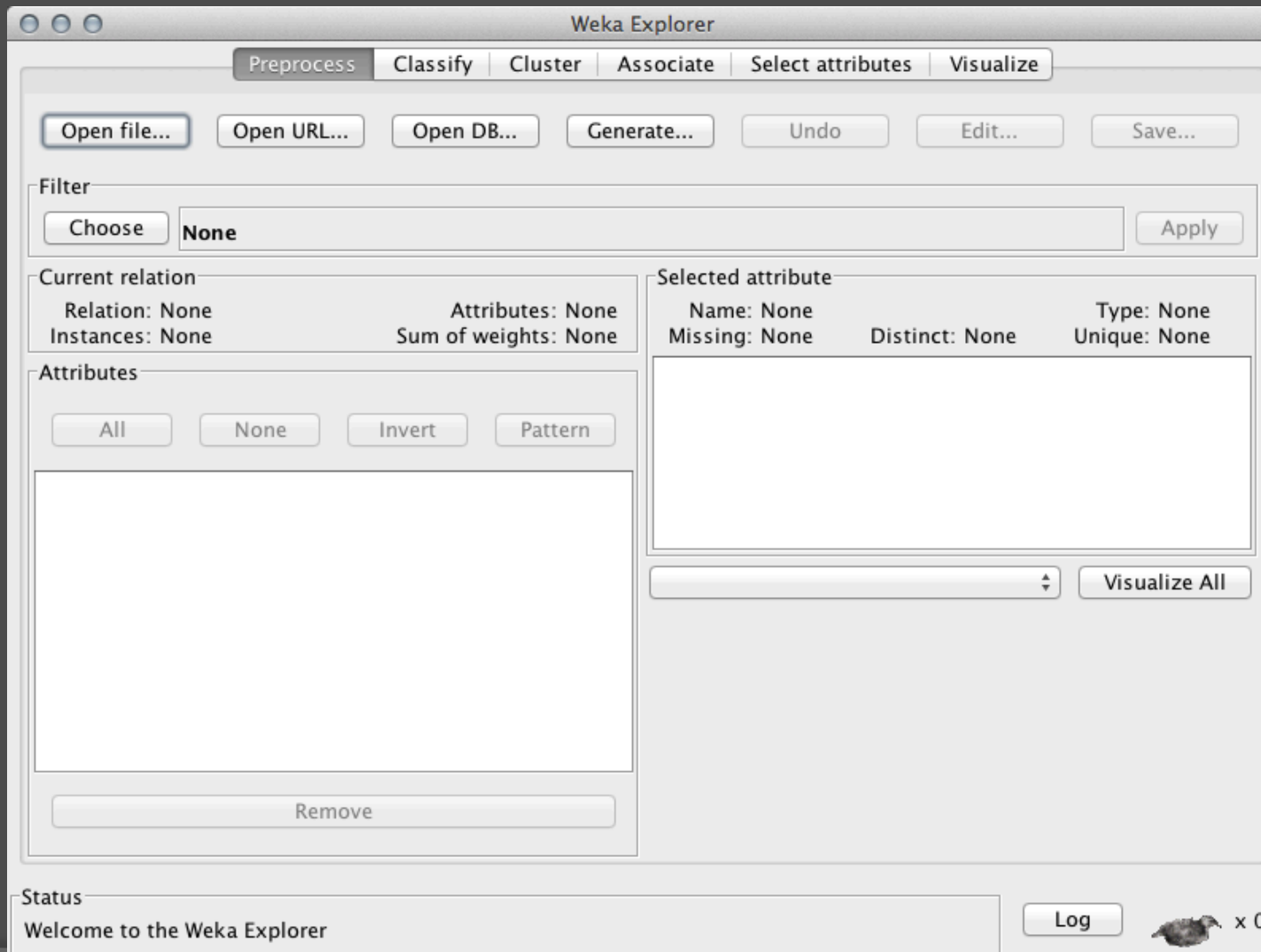
Experimenter (2)



Explorer (1)

- ⦿ Preprocess
- ⦿ Classify
- ⦿ Cluster
- ⦿ Associate
- ⦿ Select attributes
- ⦿ Visualize

Explorer (2)



Using WEKA

Data format supported by Weka

- ⦿ Data can be imported from a file in various formats:
 - ARFF – default format file
 - CSV – separated by commas or tabulations
 - C4.5 – codify under C4.5 format (.names to store the names and .data to store the data)
 - JSON – data files used by Javascript

ARFF files (1)

```
diabetes.arff
@relation pima_diabetes
@attribute 'preg' real
@attribute 'plas' real
@attribute 'pres' real
@attribute 'skin' real
@attribute 'insu' real
@attribute 'mass' real
@attribute 'pedi' real
@attribute 'age' real
@attribute 'class' { tested_negative, tested_positive}
@data
6,148,72,35,0,33.6,0.627,50,tested_positive
1,85,66,29,0,26.6,0.351,31,tested_negative
8,183,64,0,0,23.3,0.672,32,tested_positive
1,89,66,23,94,28.1,0.167,21,tested_negative
0,137,40,35,168,43.1,2.288,33,tested_positive
5,116,74,0,0,25.6,0.201,30,tested_negative
3,78,50,32,88,31,0.248,26,tested_positive
10,115,0,0,0,35.3,0.134,29,tested_negative
2,197,70,45,543,30.5,0.158,53,tested_positive
8,125,96,0,0,0,0.232,54,tested_positive
4,110,92,0,0,37.6,0.191,30,tested_negative
10,168,74,0,0,38,0.537,34,tested_positive
10,139,80,0,0,27.1,1.441,57,tested_negative
1,189,60,23,846,30.1,0.398,59,tested_positive
5,166,72,19,175,25.8,0.587,51,tested_positive
7,100,0,0,0,30,0.484,32,tested_positive
0,118,84,47,230,45.8,0.551,31,tested_positive
7,107,74,0,0,29.6,0.254,31,tested_positive
1,103,30,38,83,43.3,0.183,33,tested_negative
```

ARFF files (2)

- Header:

@relation <relation name>

- Attributes declaration:

@attribute <name> <type>

where <type> can be a value (numeric, string, date, etc) or nominal (set of values, e.g. {female, male})

- Data

@data

...

ARFF file example

% file to test.

@relation test

@attribute name STRING

@attribute health {good, bad}

@attribute weight NUMERIC

@attribute date_analyse DATE "dd-MM-yyyy HH:mm"

@data

Alice, good, 38.43, "12-04-2003 12:23"

'Maria Jose', ?, 34.53, "14-05-2003 13:45"

Alex, good, 43, "01-01-2004 08:04"

Richard, ?, ?, "03-04-2003 11:03"

ARFF files (sparse format)

- ◉ Considering only the non 0 values
- ◉ Represent each values with: POSITION VALUE information
- ◉ Each couple (POSITION - VALUE) is separated with a comma
- ◉ Useful for documents representation in ARFF format

e.g.

@data

0, X, Y, "male"

0, 0, W, "male"

→

→

@data

{1 X, 2 Y, 3 "male"}

{2 W, 3 "male"}

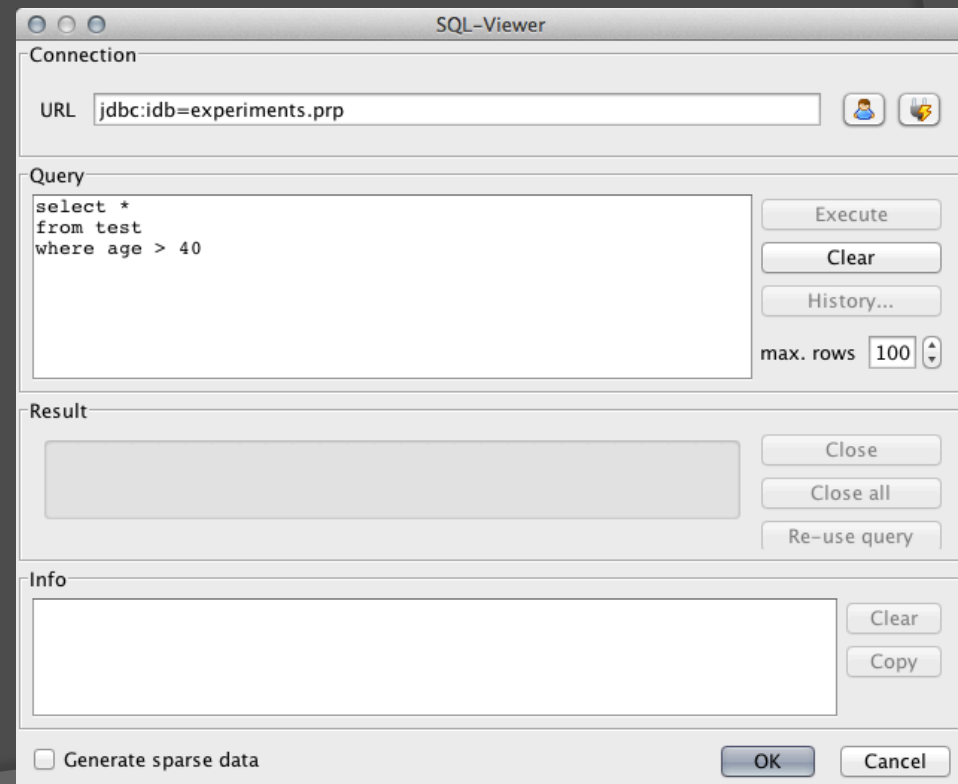
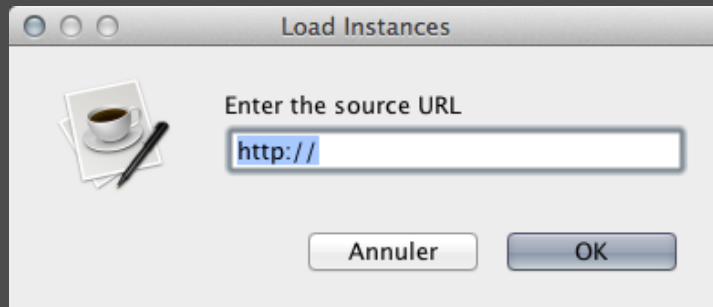
Exercise

- Represent the following table in ARFF file (simple and sparse formats)

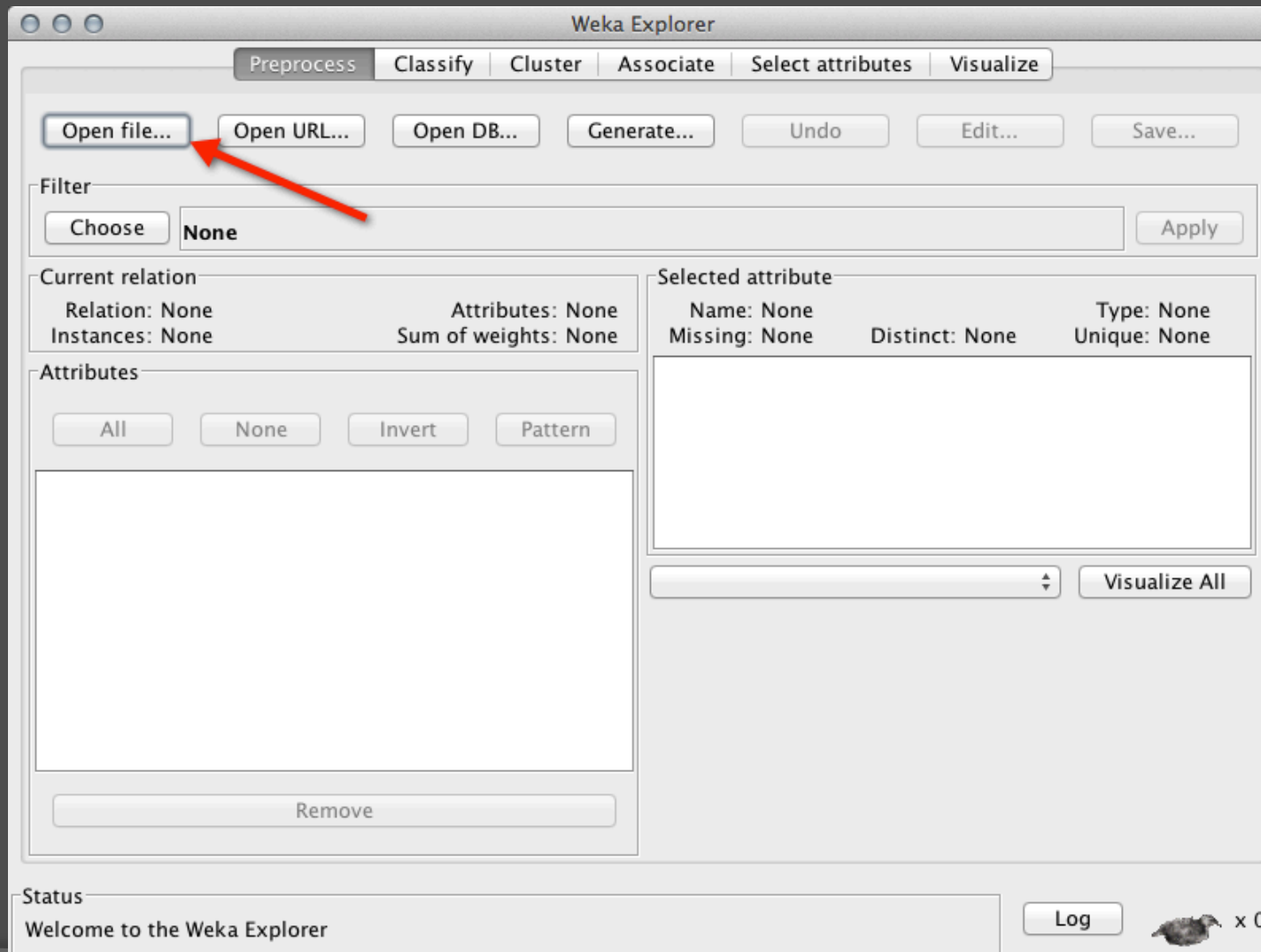
City	Date	Temperature	Humidity	Wind	Emergency
Alès	03/14	14.4	68	57	Yes
Paris	03/15	18.4	60		No
Nîmes	03/14	20.3	72	45	Yes
Nice	04/01	15.6	68	11	No
Lunel	03/18	28.0	71		No

Open a DB or URL

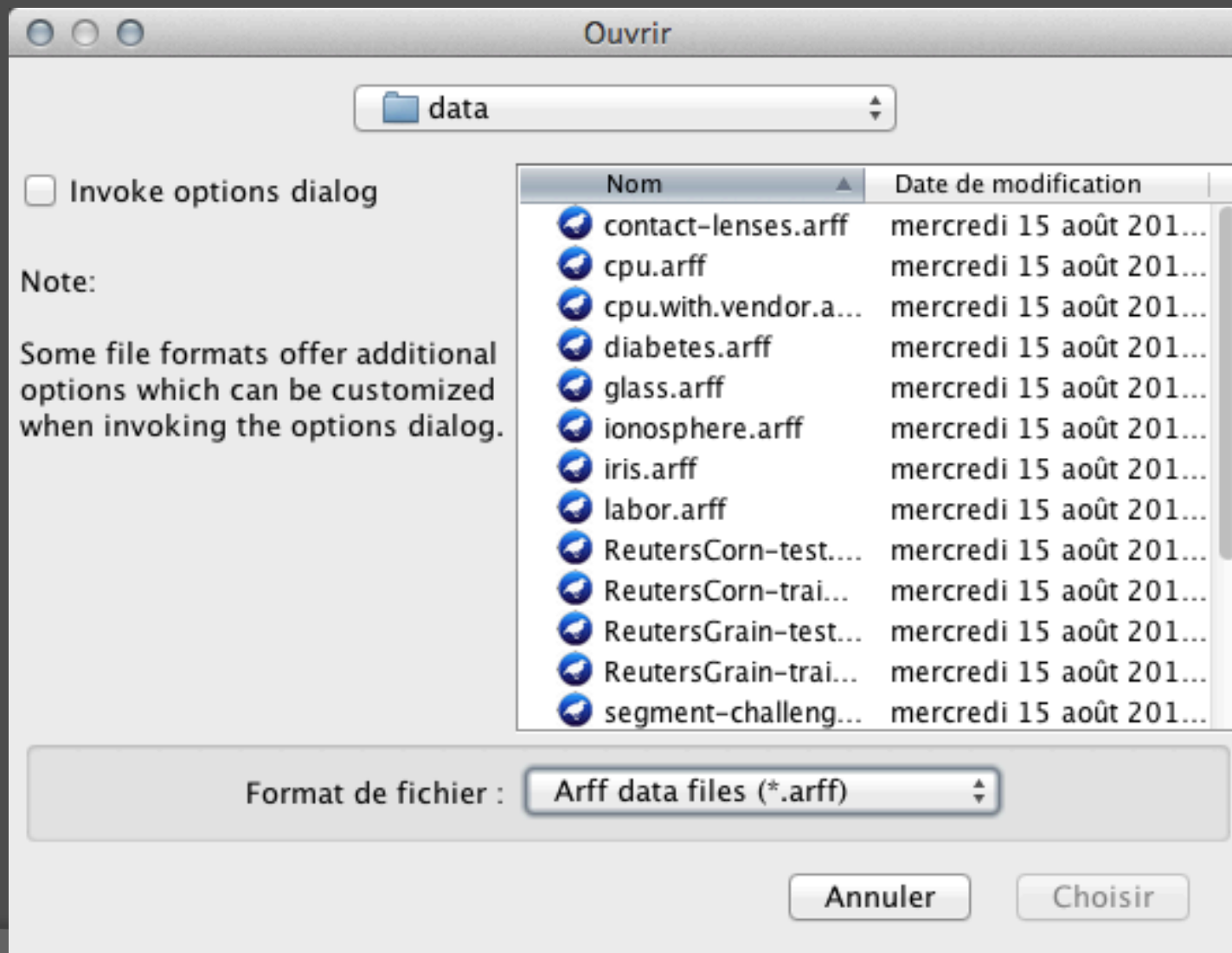
- Data can also be read from a URL or from an SQL database (using JDBC)



Open file using Explorer (1)



Open file using Explorer (2)



Open file using Explorer (3)

The screenshot shows the Weka Explorer application window. The 'Preprocess' tab is active. The 'Current relation' is 'weather' with 14 instances and 5 attributes. The 'Attributes' list includes 'outlook', 'temperature', 'humidity', 'windy', and 'play'. The 'Selected attribute' is 'outlook', which is a nominal attribute with 3 distinct values and 5 instances. The 'Class' is 'play (Nom)'. A stacked bar chart visualizes the data for the 'outlook' attribute, showing the distribution of 'play' values for each 'outlook' category.

Weka Explorer

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

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

Filter: Choose None Apply

Current relation: Relation: weather, Instances: 14, Attributes: 5, Sum of weights: 14

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

No.	Label	Count	Weight
1	sunny	5	5.0
2	overcast	4	4.0
3	rainy	5	5.0

Attributes: All None Invert Pattern

No.	Name
<input checked="" type="checkbox"/>	1 outlook
<input type="checkbox"/>	2 temperature
<input type="checkbox"/>	3 humidity
<input type="checkbox"/>	4 windy
<input type="checkbox"/>	5 play

Remove

Class: play (Nom) Visualize All

5 4 5

Status: OK Log x 0

Open file using Explorer (4)

The screenshot shows the Weka Explorer interface with the 'weather' dataset loaded. The 'Attributes' list on the left includes 'outlook', 'temperature', 'humidity', 'windy', and 'play'. The 'Classes' dropdown is set to 'play (Nom)'. The 'Attribute characteristics' table for 'outlook' shows 3 distinct values: 'sunny' (5 instances), 'overcast' (4 instances), and 'rainy' (5 instances). Three red arrows point to the 'Attributes' list, the 'Classes' dropdown, and the 'Attribute characteristics' table.

Attributes

No.	Name
1	<input checked="" type="checkbox"/> outlook
2	<input type="checkbox"/> temperature
3	<input type="checkbox"/> humidity
4	<input type="checkbox"/> windy
5	<input type="checkbox"/> play

Classes

Class: play (Nom)

Attribute characteristics

No.	Label	Count	Weight
1	sunny	5	5.0
2	overcast	4	4.0
3	rainy	5	5.0

Current relation

Relation: weather
Instances: 14
Attributes: 5
Sum of weights: 14

Status

OK

Log x 0

Pre-processing tools (1)

- Pre-processing tools in WEKA are called “filters”
- WEKA contains filters for:

Discretization, normalization, resampling, attribute selection, transforming and combining attributes, ...

Pre-processing tools (2)

The screenshot shows the Weka Explorer interface in the Preprocess tab. The 'Filter' section is set to 'None', indicated by a red arrow. The 'Current relation' is 'weather' with 14 instances and 5 attributes. The 'Attributes' list shows 'outlook' selected. The 'Selected attribute' section displays a table for 'outlook' with 3 distinct values: sunny (5 instances), overcast (4 instances), and rainy (5 instances). Below this, a bar chart visualizes the distribution of the 'play' class across the 'outlook' categories.

Current relation
Relation: weather
Instances: 14
Attributes: 5
Sum of weights: 14

Attributes
All None Invert Pattern

No.	Name
<input checked="" type="checkbox"/>	1 outlook
<input type="checkbox"/>	2 temperature
<input type="checkbox"/>	3 humidity
<input type="checkbox"/>	4 windy
<input type="checkbox"/>	5 play

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

No.	Label	Count	Weight
1	sunny	5	5.0
2	overcast	4	4.0
3	rainy	5	5.0

Class: play (Nom) Visualize All

Status: OK Log x 0

Pre-processing tools (3)

The screenshot shows the Weka Explorer application window. The 'Preprocess' tab is active, displaying a list of pre-processing filters. A filter selection dialog is open, showing a tree view of filters under the 'weka' directory. The 'attribute' sub-directory is expanded, showing various filters like 'AddClassification', 'AttributeSelection', 'ClassOrder', 'Discretize', 'NominalToBinary', 'Resample', 'SpreadSubsample', and 'StratifiedRemoveFolds'. The 'instance' sub-directory contains 'Resample', 'SpreadSubsample', and 'StratifiedRemoveFolds'. The 'unsupervised' sub-directory is also expanded, showing 'attribute' filters like 'Add', 'AddCluster', 'AddExpression', and 'AddID'. The 'AttributeSelection' filter is highlighted in the dialog.

The main window shows the 'Selected attribute' section for 'outlook'. The attribute is of type 'Nominal' with 3 distinct values and 0 missing values. The data is summarized in the following table:

No.	Label	Count	Weight
1	sunny	5	5.0
2	overcast	4	4.0
3	rainy	5	5.0

The 'Class: play (Nom)' dropdown is set to 'play (Nom)'. Below the table, there are three stacked bar charts representing the distribution of the 'play' class for each 'outlook' value. The first chart (sunny) has a total height of 5, with a red top section of 5 and a blue bottom section of 0. The second chart (overcast) has a total height of 4, with a red top section of 0 and a blue bottom section of 4. The third chart (rainy) has a total height of 5, with a red top section of 5 and a blue bottom section of 0.

Buttons at the bottom include 'Filter...', 'Remove filter', 'Close', 'Log', and a small animal icon with 'x 0'.

Example – Discretization (1)

- To obtain categorical data
- Used on numerical attributes
 1. Open a file (weather.arff for example)
 2. Choose a filter : *Filters* → *unsupervised* → *discretize* **1**
 3. Left-click on properties **2**
 4. Change the number of *binds* and *useEqualFrequency*
 5. Click on OK and APPLY **3**

Example – Discretization (2)

The screenshot shows the Weka Explorer interface with the 'Discretize' filter applied to the 'outlook' attribute. The filter parameters are set to '-F -B 4 -M -1.0 -R first-last'. The 'outlook' attribute is selected, and its distribution is shown in a table and a bar chart.

Filter: Discretize -F -B 4 -M -1.0 -R first-last

Current relation: Relation: weather-weka.filters... Attributes: 5 Instances: 14 Sum of weights: 14

Attributes: outlook, temperature, humidity, windy, play

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

No.	Label	Count	Weight
1	sunny	5	5.0
2	overcast	4	4.0
3	rainy	5	5.0

Class: play (Nom)

Bar Chart: The bar chart shows the distribution of the 'play' class for each 'outlook' category. The 'sunny' category has 5 instances (red bar), the 'overcast' category has 4 instances (blue bar), and the 'rainy' category has 5 instances (red bar).

Status: OK

Exercise

In the example, compare the characteristics of attributes before and after discretization:

- Evaluate the *outlook* attribute characteristics: Comment the results
- Evaluate the *temperature* attribute characteristics: Comment the results
- Export the results into *arff* and *csv* files

Data normalization

- ⦿ Pre-processing technique
- ⦿ The filter *standardization* allow us standardize all numerical values of the data set into values belonging the interval $[0, 1]$
- ⦿ For more information, see “More”

Example (1)

1. Open *fruitfly.arff* 1
2. See the dataset using the button *Edit* 2
3. Choose: *filters* → *unsupervised* → *attribute* → *normalize* 3
4. Set *scale* to 1.0 into options 4
5. Click on *Apply* 5
6. See the data using the button *Edit* 6

Example (2)

The screenshot shows the Weka Explorer application window. The title bar reads "Weka Explorer". The menu bar includes "Preprocess", "Classify", "Cluster", "Associate", "Select attributes", and "Visualize". The "Preprocess" menu is currently selected, showing options: "Open file...", "Open URL...", "Open DB...", "Generate...", "Undo", "Edit...", and "Save...".

Below the menu bar is the "Filter" section, which contains a "Choose" button, a text field with the filter expression "Normalize -S 1.0 -T 0.0", and an "Apply" button. Red arrows labeled "2,6" point to the "Apply" button and the "Edit..." button in the menu bar.

The "Current relation" section displays "Relation: fruitfly", "Instances: 125", "Attributes: 5", and "Sum of weights: 125".

The "Attributes" section has buttons for "All", "None", "Invert", and "Pattern". Below these is a list of attributes:

No.	Name
1	<input checked="" type="checkbox"/> PARTNERS
2	<input type="checkbox"/> TYPE
3	<input type="checkbox"/> THORAX
4	<input type="checkbox"/> SLEEP
5	<input type="checkbox"/> class

A red arrow labeled "3" points to the "PARTNERS" attribute in this list. A "Remove" button is located below the list.

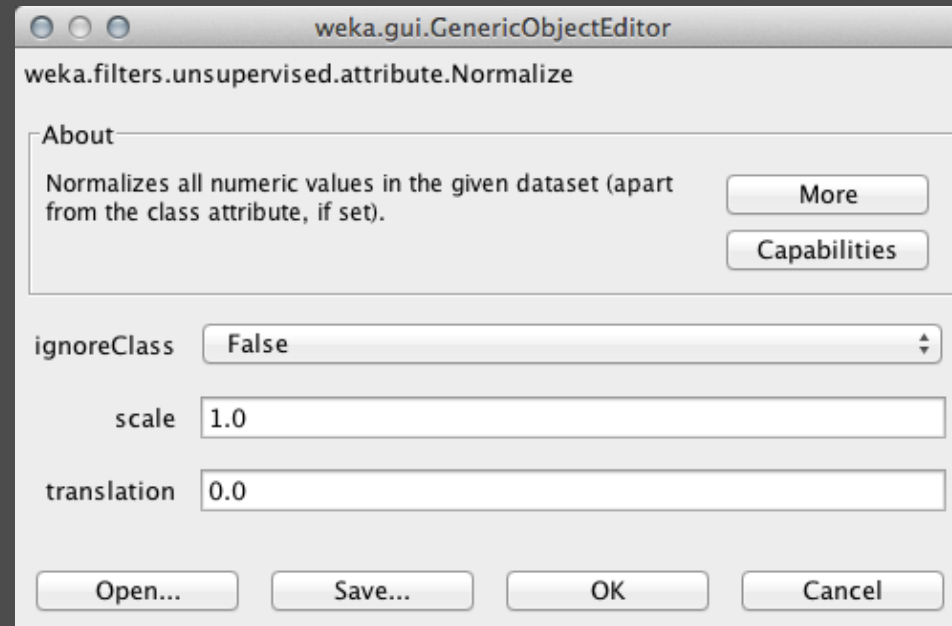
The "Selected attribute" section shows "Name: PARTNERS", "Missing: 0 (0%)", "Distinct: 3", and "Type: Nominal Unique: 0 (0%)". Below this is a table:

No.	Label	Count	Weight
1	8	50	50.0
2	0	25	25.0
3	1	50	50.0

A red arrow labeled "5" points to the "Weight" column of this table. Below the table is a "Class: class (Num)" dropdown menu and a "Visualize All" button.

At the bottom of the window is a "Status" bar showing "OK" and a "Log" button. A small bird icon and "x 0" are also visible in the bottom right corner.

Example (3)



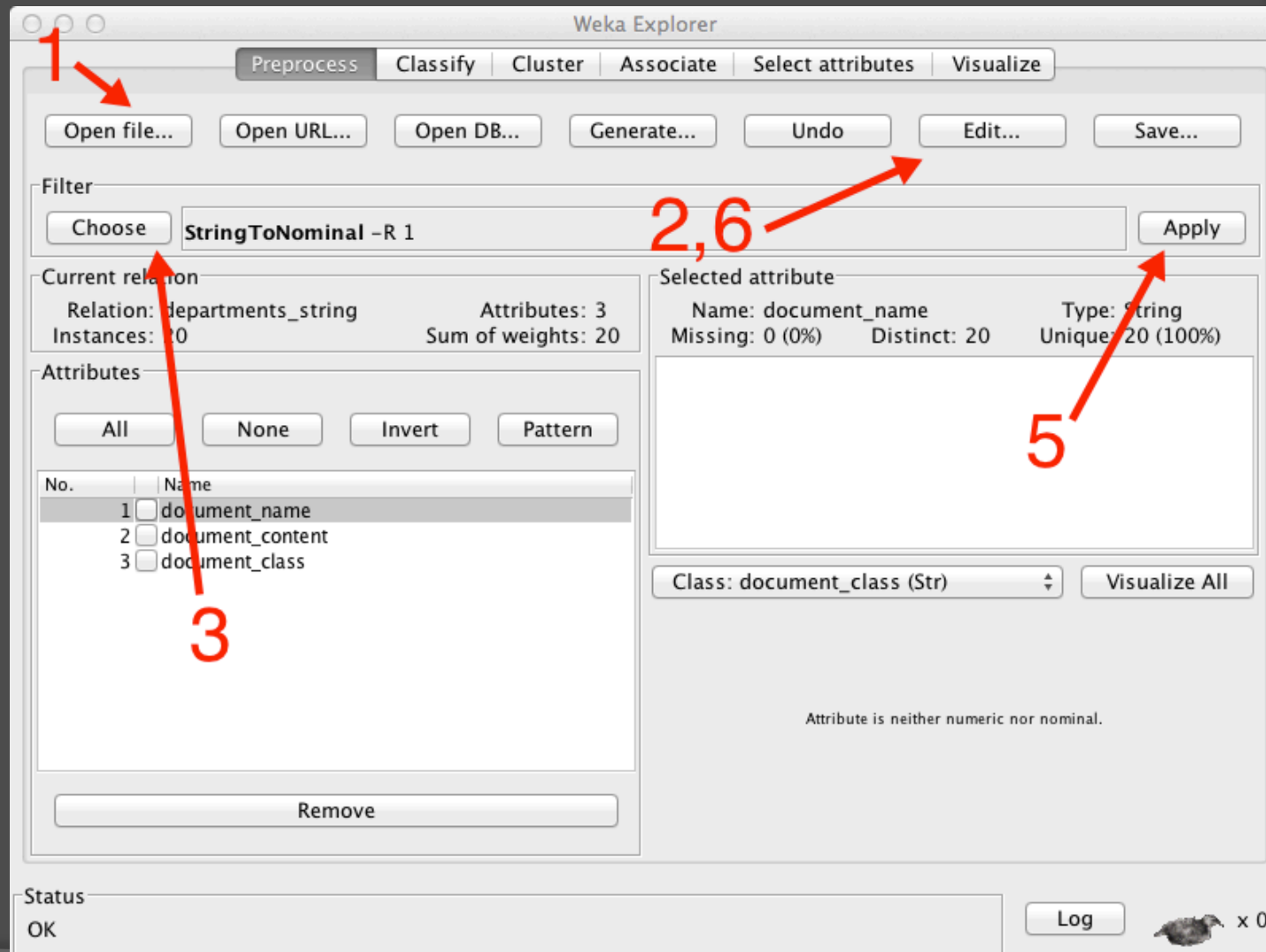
String attribute into nominal

- ⦿ Pre-processing tool
- ⦿ Converting a string attribute into nominal
- ⦿ Finite number of values (string)
- ⦿ For more information, see “More”

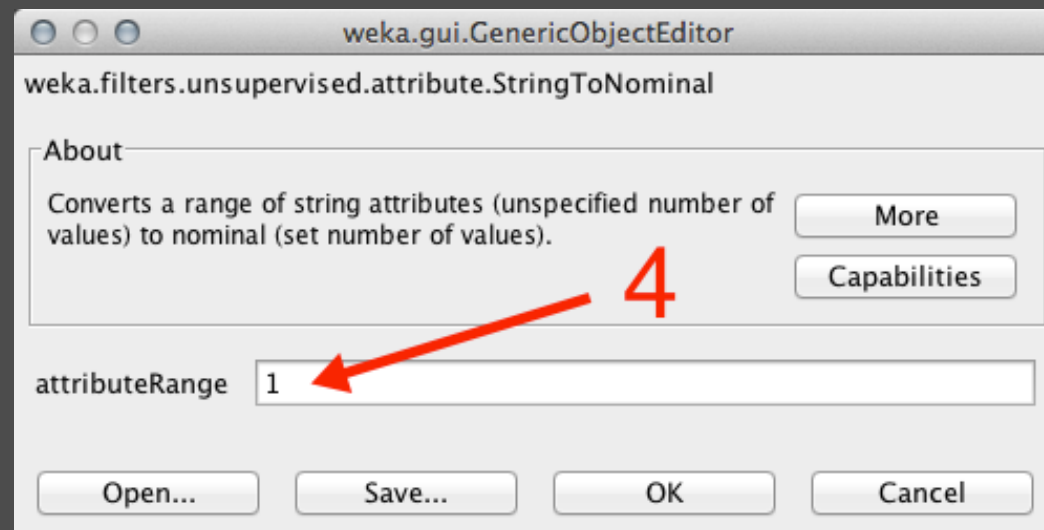
Example (1)

1. Open *Departments-string.arff* **1**
2. See the dataset using the button *Edit* **2**
3. Choose: *filters* → *unsupervised* → *attribute* → *StringToNominal* **3**
4. Set *attributeRange* to 1 into options **4**
5. Click on *Apply* **5**
6. See the dataset using the button *Edit* **6**

Example (2)



Example (3)



Text data to String vector (1)

- Pre-processing technique
- Converting text data into TF-IDF (Term Frequency – Inverted Document Frequency) attribute format
- Used on string data
- For more information, see “More”

Example (1)

1. Open *Departments-string.arff* **1**
2. See the dataset using the button *Edit* **2**
3. Choose: *filters* → *unsupervised* → *attribute* → *StringToWordVector* **3**
4. Set some options **4**
5. Click on *Apply* **5**
6. See the dataset using the button *Edit* **6**

Example (2)

The screenshot shows the Weka Explorer application window. The interface includes a menu bar with options: Preprocess, Classify, Cluster, Associate, Select attributes, and Visualize. Below the menu bar are buttons for Open file..., Open URL..., Open DB..., Generate..., Undo, Edit..., and Save....

The Filter section contains a Choose button, a text field with the filter `StringToWordVector -R first-last -W 1000 -prune-rate -1.0 -N 0 -stemmer weka.core.stemmer`, and an Apply button.

The Current relation section shows: Relation: departments_string, Instances: 20, Attributes: 3, Sum of weights: 20.

The Attributes section has buttons for All, None, Invert, and Pattern. Below is a table:

No.	Name
1	<input checked="" type="checkbox"/> document_name
2	<input type="checkbox"/> document_content
3	<input type="checkbox"/> document_class

Below the table is a Remove button.

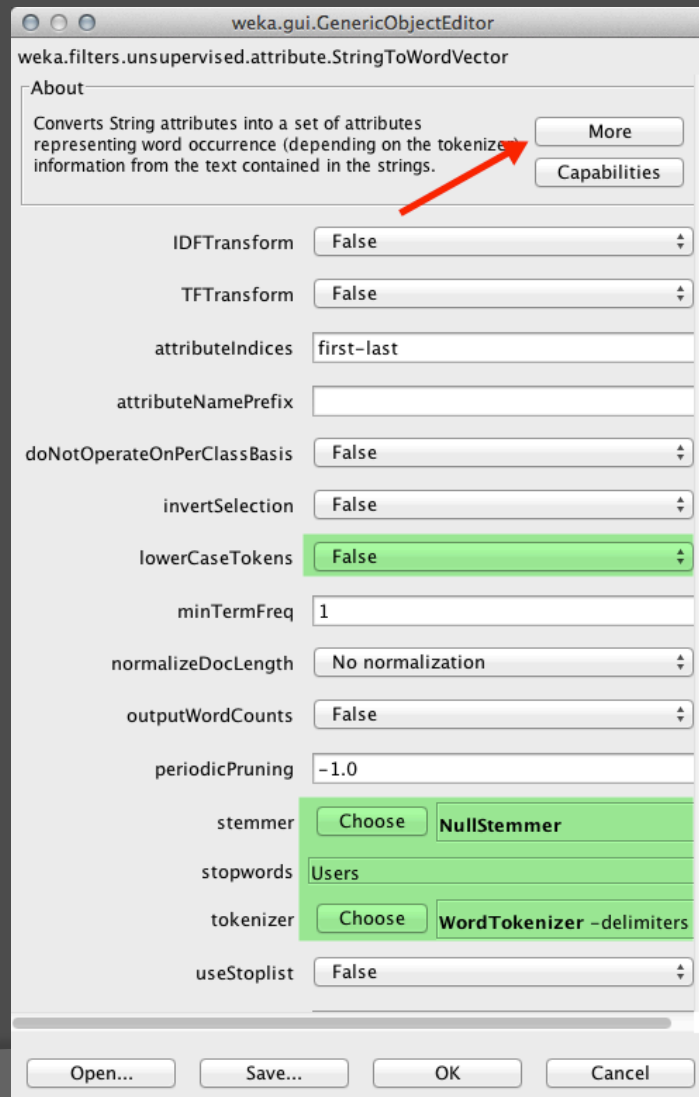
The Selected attribute section shows: Name: document_name, Type: String, Missing: 0 (0%), Distinct: 20, Unique: 20 (100%).

Below this section is a dropdown menu set to 'Class: document_class (Str)' and a Visualize All button.

At the bottom, the Status bar shows 'OK' and a Log button. A small dog icon and 'x 0' are also visible.

Red annotations are present: '1' points to the window title bar; '2,6' points to the filter text field; '3' points to the 'document_name' row in the attributes table; '5' points to the 'Apply' button; and '6' points to the 'Edit...' button.

Example (3)



Example (4)

```
@relation departments_string

@attribute document_name string
@attribute document_content string
@attribute document_class string

@data

Anthropology, "anthropology anthropology anthropology co
archaeology and linguistics beyond these subfields conce
comparison the anthropology major provides students with
a range of careers from public service to marketing and
with faculty doing research students regularly attend pr
special programs include summer field schools in archaeo
sponsored diversity training institutes program of study
department website", A
Art, "art art the art department s undergraduate degree
printmaking sculpture ceramics and graphic design both c
encourages all forms of creative explorations via a deve
successful portfolio review of at least 10 works to beco
developed on an individual basis consistent with the goa
the samuel chen art center a gallery that offers regular
known artists sol lewitt cleve gray and robert cottingha
museums and galleries programs of study ba ms department
department website", B
Biology, "biology biological sciences the undergraduate
explore the discipline broadly specialized undergraduate
interpretation also available are specialized graduate p
various health and medical professions are advised prima
programs require a research project or internship many l
environmental rooms cell culture and protein purificatio
computer laboratory are available for research and instr
copernicus hall 332 phone 832 2645 department website fu
Chemistry. "chemistry chemistry the chemistry department
```



```
@attribute welte numeric
@attribute western numeric
@attribute willard numeric
@attribute wolff numeric
@attribute women numeric
@attribute works numeric
@attribute writing numeric
@attribute year numeric
@attribute york numeric

@data
{2 1,10 1,25 1,27 1,28 1,39 1,62 1,66 1,67 1,70 1,71 1
1,147 1,153 1,159 1,160 1,162 1,164 1,183 1,192 1,193
1,286 1,291 1,297 1,303 1,306 1,320 1,321 1,333 1,338
1,393 1,403 1,404 1,411 1,414 1,420 1,425 1,431 1,433
{26 1,39 1,52 1,56 1,61 1,62 1,71 1,85 1,96 1,104 1,11
1,283 1,288 1,310 1,320 1,322 1,324 1,328 1,346 1,366
1,478 1,485 1,486 1,487 1,491 1,494 1,495 1,509 1,510
1,558 1,561 1,566 1,568 1,571 1,578 1,579 1,582 1,587
1,629 1,630 1,631 1,632 1,636 1,646 1,653 1,654 1,661
1,725 1,734 1,735 1,739 1,754 1,755 1,756 1,762 1}
{11 1,21 1,26 1,27 1,29 1,39 1,51 1,52 1,56 1,58 1,62
1,141 1,148 1,152 1,153 1,160 1,181 1,186 1,188 1,191
1,258 1,260 1,269 1,270 1,283 1,284 1,292 1,302 1,310
1,375 1,380 1,392 1,393 1,395 1,396 1,398 1,399 1,404
{12 1,23 1,26 1,27 1,30 1,60 1,69 1,73 1,86 1,89 1,94
1,151 1,153 1,172 1,177 1,181 1,214 1,218 1,219 1,225
1,366 1,368 1,382 1,399 1,402 1,404 1,408 1,413 1,415
{0 1,26 1,39 1,41 1,52 1,60 1,61 1,62 1,72 1,85 1,86 1
1,218 1,228 1,240 1,248 1,258 1,283 1,309 1,310 1,320
1,399 1,420 1,421 1,431 1,433 1,435 1,444 1,445 1,446
1,547 1,554 1,562 1,577 1,581 1,586 1,588 1,590 1,597
1,704 1,719 1,721 1,723 1,741 1,744 1,745 1,752 1}
```

Attribute selection (1)

- ① The most useful part of this is attribute selection (also called feature selection)
- ② Select relevant attributes
- ③ Remove redundant and/or irrelevant attributes

Attribute selection (1)

Objectives:

- ① Simpler model
 - More transparent
 - Easier to interpret
- ② Faster model induction
- ③ Structural knowledge
 - Knowing which attributes are important may be inherently important to the application

Attribute selection (2)

Attribute Evaluator

Search Method	Attribute Evaluator	
	Attributes	Subsets of attributes
Best First		YES
Greedy Stepwise		YES
Ranker	YES	

Attribute selection (3)

Filters:

- ⦿ Ranked list of attributes
 - Typical when each attribute is evaluated individually
- ⦿ A selected subset of attributes
 - Greedy Stepwise and Best first
 - Random search such as genetic algorithm

Example (1)

1. Open *diabetes.arff* 1
2. Choose: *filters* → *supervised* → *attribute*
→ *AttributeSelection* 2

It's possible to use “*SelectAttributes*” tab

3. Set some options 3
4. Click on *Apply* 4

Example (2)

Weka Explorer

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

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

Filter: Choose **AttributeSelection** -E "weka.attributeSelection.CfsSubsetEval" -S "weka.attributeSelection.Be: Apply

Current Relation: Relation: pima_diabetes-wek... Attributes: 5 Sum of weights: 768
Instances: 768

Attributes: All | None | **Invert** | Pattern

No.	Name
<input checked="" type="checkbox"/> 1	plas
<input type="checkbox"/> 2	mass
<input type="checkbox"/> 3	pedi
<input type="checkbox"/> 4	age
<input type="checkbox"/> 5	class

Remove

Selected attribute: Name: plas Type: Numeric
Missing: 0 (0%) Distinct: 136 Unique: 19 (2%)

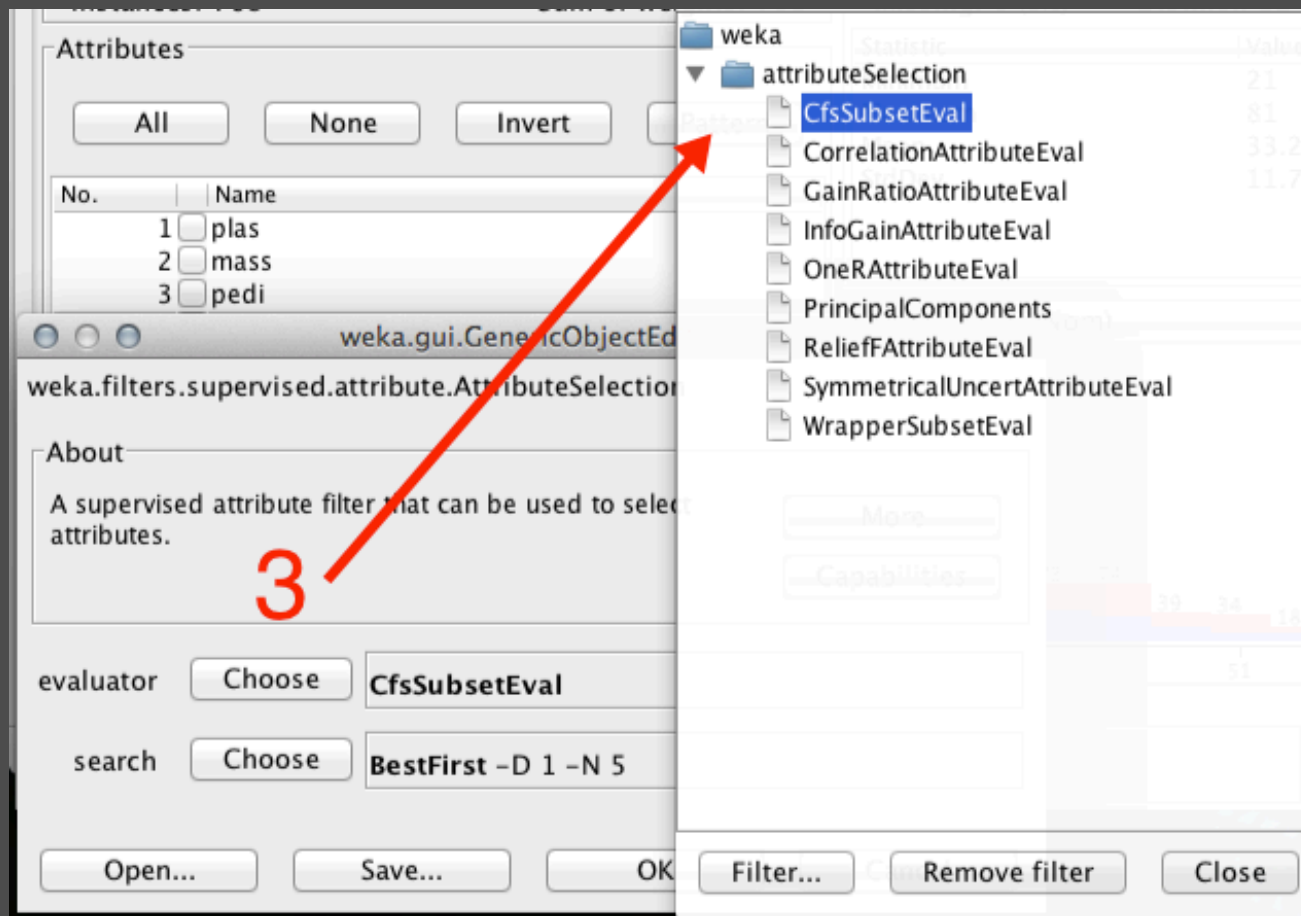
Statistic	Value
Minimum	0
Maximum	199
Mean	120.89
StdDev	31.973

Class: class (Nom) Visualize All

Status: OK Log x 0

Red arrows and numbers 1, 2, and 4 highlight the 'Invert' button, the 'plas' attribute in the list, and the 'Apply' button, respectively.

Example (3)



Exercise

- ① Discuss the results after applying *AttributeSelection* on *diabetes.arff* using “default” parameters
- ① Change some parameters and compare the results
- ① Apply a classification algorithm (e.g., J48) aux datasets with/without attribute selection. Compare the results

Conclusion

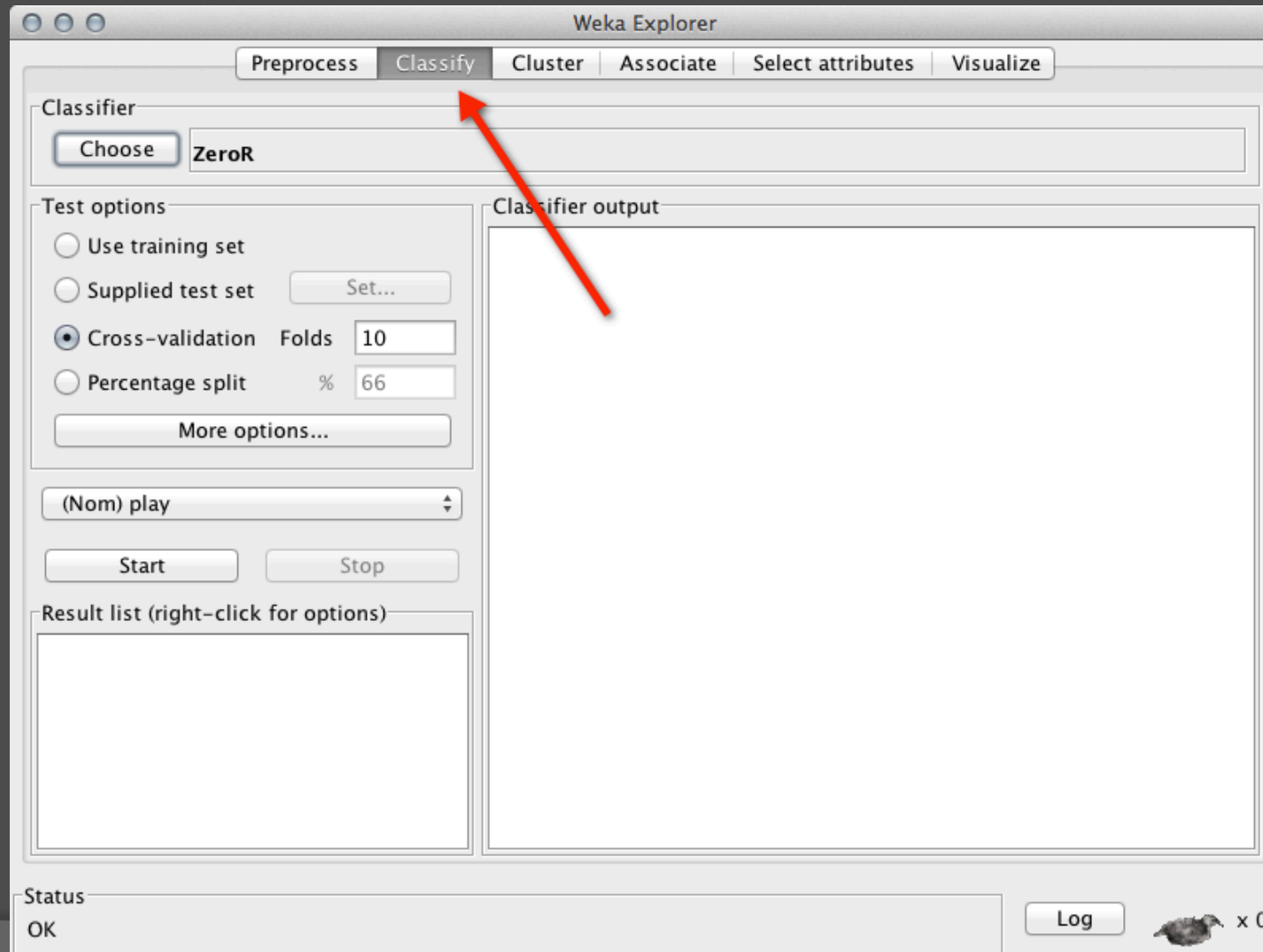
Data preprocessing is very important, and it has an important impact on the quality of learning process

Classifiers

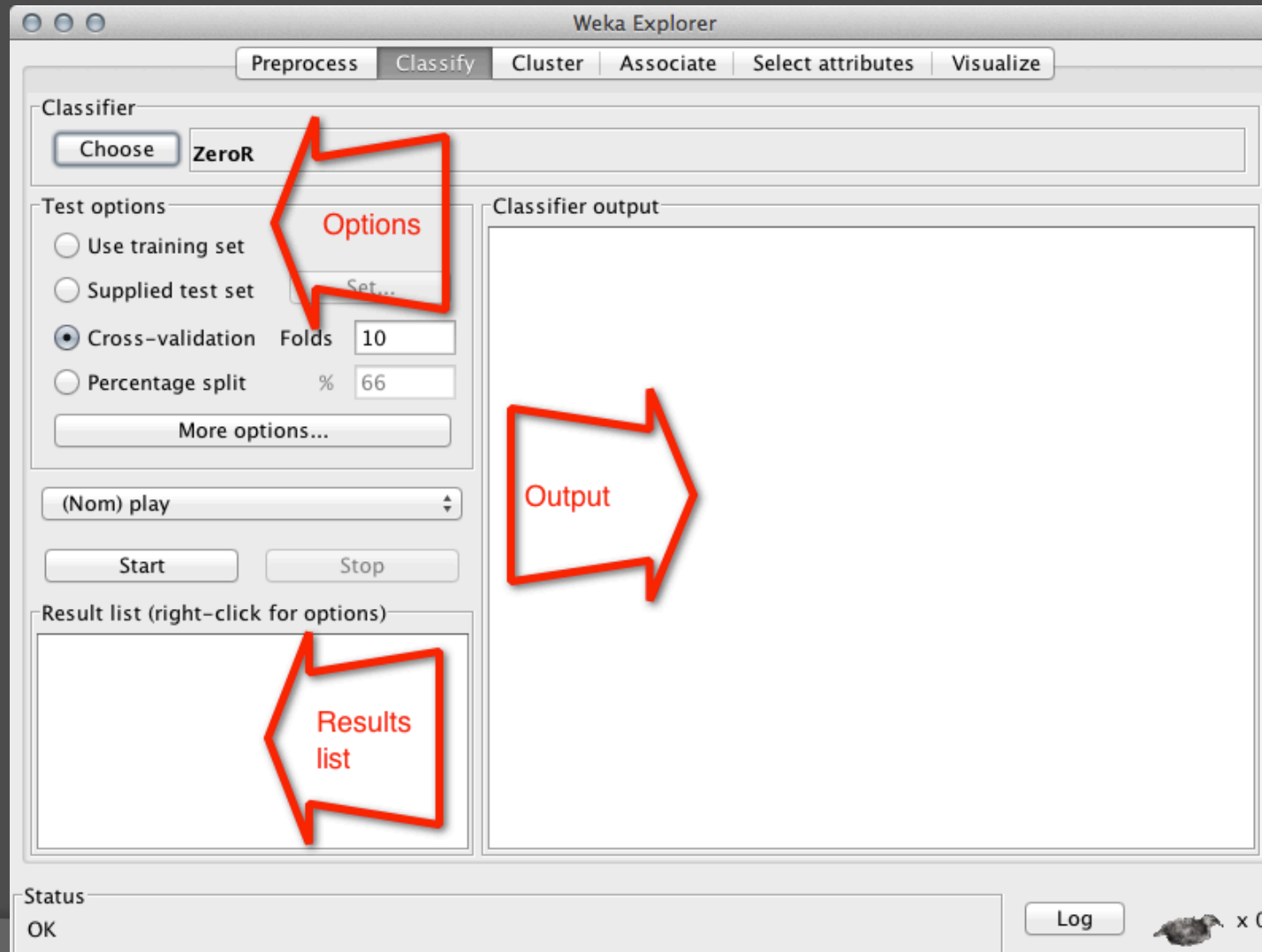
Classifiers in Weka (1)

- ⦿ Classifiers in WEKA are models for predicting nominal or numeric quantities
- ⦿ Classification algorithms include:
 - Decision trees
 - Naïve Bayes Classification
 - Support vector machine (SVM)
 - Multi-layer perceptron
 - Bayes network, etc.
- ⦿ Meta-classifiers:
 - Combination
 - Bagging
 - Boosting, etc.

Classifiers in Weka (2)



Classifiers in Weka (3)



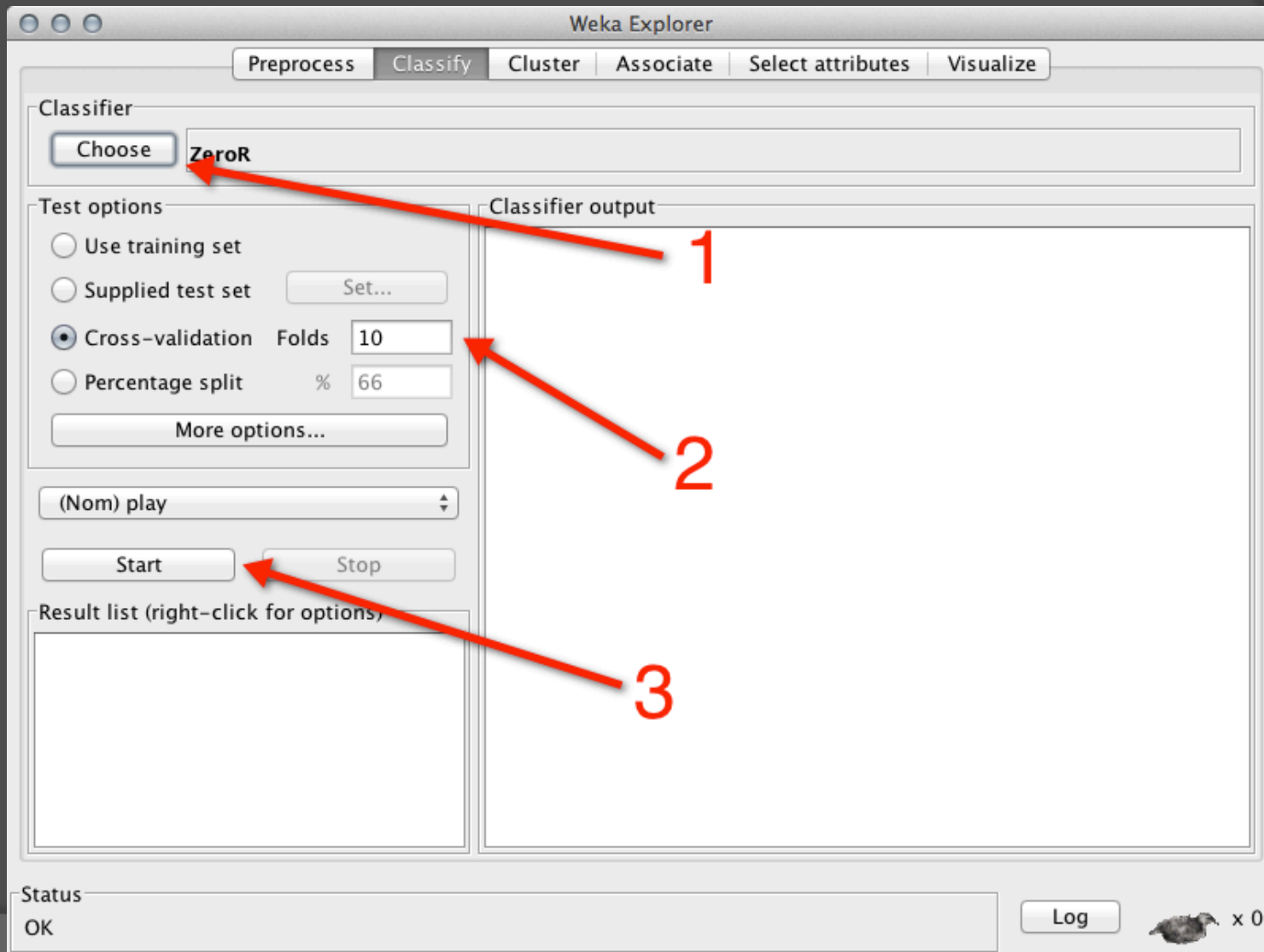
Classifiers : options

- ④ **Training set:** the classifier is evaluated on how well it predicts the class of the instances it was trained on
- ④ **Supplied test set:** the classifier is evaluated on how well it predicts the class of a set of instances loaded from a file
- ④ **Cross-Validation:** the classifier is evaluated by cross-validation, using the number of folds that are entered in the Folds text field
- ④ **Percentage Split:** the classifier is evaluated on how well it predicts a certain percentage of the data which is held out for testing

Example (1)

1. Open iris.arff
2. Go to *Classify* tab
3. Choose a classifier : *Classifier* → *Bayes* → *NaiveBayes* 1
4. Set *Cross-Validation* value to 10 2
5. Click on *Start* button 3

Example (2)



Example (3)

The screenshot shows the Weka Explorer interface with the 'Classify' tab selected. The classifier chosen is 'NaiveBayes'. The test options are set to 'Cross-validation' with 10 folds and a 70% split. The classifier output shows the following information:

```
=== Run information ===
Scheme:      weka.classifiers.bayes.NaiveBayes
Relation:    iris
Instances:   150
Attributes:  5
              sepallength
              sepalwidth
              petallength
              petalwidth
              class
Test mode:   10-fold cross-validation

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

Naive Bayes Classifier

Attribute          Class
                   Iris-setosa Iris-versicolor Iris-virginica
                   (0.33)      (0.33)      (0.33)
-----
sepallength
  mean              4.9913      5.9379      6.5795
  std. dev.         0.355       0.5042      0.6353
  weight sum        50          50          50
  precision         0.1059     0.1059     0.1059
sepalwidth
  mean              3.4015     2.7687     2.9629
  std. dev.         0.3925     0.3038     0.3088
  weight sum        50          50          50
  precision         0.1091     0.1091     0.1091
```

The result list on the left shows the following entries:

- 15:39:43 - rules.OneR
- 15:42:00 - rules.OneR
- 15:48:21 - trees.J48
- 15:52:07 - trees.J48
- 15:53:03 - trees.J48
- 16:03:38 - bayes.NaiveBayes
- 16:33:27 - trees.DecisionStump
- 16:35:20 - bayes.NaiveBayes

The status bar at the bottom shows 'OK' and a 'Log' button.

Interpretation of results (1)

```
=== Evaluation on test split ===  
Time taken to test model on training split: 0.04 seconds  
  
=== Summary ===  
Correctly Classified Instances      43      95.5556 %  
Incorrectly Classified Instances    2      4.4444 %  
Kappa statistic                    0.9331  
Mean absolute error                0.0375  
Root mean squared error            0.158  
Relative absolute error             8.422 %  
Root relative squared error        33.4987 %  
Coverage of cases (0.95 level)     97.7778 %  
Mean rel. region size (0.95 level) 37.037 %  
Total Number of Instances          45  
  
=== Detailed Accuracy By Class ===  
      TP Rate  FP Rate  Precision  Recall  F-Measure  MCC   ROC Area  PRC Area  Class  
      1      0      1          1      1          1     1          1          1  Iris-setosa  
      1      0.069  0.889     1      0.941     0.91  0.987     0.976     1  Iris-versicolor  
      0.867  0      1          0.867  0.929     0.901 0.987     0.979     1  Iris-virginica  
Weighted Avg.  0.956  0.025    0.96     0.956    0.955    0.935 0.991     0.984  
  
=== Confusion Matrix ===  
  a  b  c  <-- classified as  
14  0  0  |  a = Iris-setosa  
 0 16  0  |  b = Iris-versicolor  
 0  2 13  |  c = Iris-virginica
```

1

2

3

Interpretation of results (2)

=== Summary === 1

- ⦿ This gives the error levels when applying the classifier.
- ⦿ The most important figures here are the numbers of correctly and incorrectly classified instances.
- ⦿ With the exception of the Kappa statistic, the remaining statistics compute various error measures based on the class probabilities assigned by the tree.

Interpretation of results (3)

=== Detailed Accuracy By Class === 2

- ⦿ The percentage of correctly classified instances is often called accuracy or sample accuracy.
- ⦿ Accuracy has some disadvantages as a performance estimate (not chance corrected, not sensitive to class distribution)
- ⦿ Area under the ROC curve is an interesting measure.

Interpretation of results (4)

=== Confusion Matrix === 3

- ⦿ This shows for each class, how instances from that class received the various classifications.
- ⦿ a, b and c representing the class labels. Here there were 45 instances, so the percentages and raw numbers add up, $aa+bb+cc = 43$, $ab+ba+ac+ca+\dots = 2$.

Decision trees

- ⦿ Learning by partitioning
- ⦿ We want to build homogeneous subgroups in terms of a nominal variable to be predicted (target) using a set of discriminant variables
- ⦿ Result must be readable
- ⦿ It must be able to automatically select discriminating variables

Example (1)

1. Open iris.arff
2. Go to *Classify* tab
3. Choose a classifier : *Classifier* → *trees* → *J48* 1
4. Set *Cross-Validation* value to 10 2
5. Click on *Start* button 3

Example (2)

The screenshot shows the Weka Explorer application window with the 'Classify' tab selected. The interface includes several sections:

- Classifier:** A dropdown menu showing 'J48 -C 0.25 -M 2'.
- Test options:** Radio buttons for 'Use training set', 'Supplied test set', 'Cross-validation', and 'Percentage split'. The 'Cross-validation' option is selected with 'Folds' set to 10.
- Classifier output:** A text area displaying run information and a decision tree model. The run information includes: Scheme: weka.classifiers.trees.J48 -C 0.25 -M 2, Relation: iris, Instances: 150, Attributes: 5 (sepalength, sepalwidth, petallength, petalwidth, class). The test mode is '10-fold cross-validation'. The classifier model is a 'J48 pruned tree' with the following structure:

```
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)
```

Additional statistics shown are 'Number of Leaves : 5' and 'Size of the tree : 9'.
- Result list (right-click for options):** A list box containing one entry: '17:18:38 - trees.J48'.

Four red arrows with numbers 1, 2, 3, and 4 point to specific elements: 1 points to the 'Classify' tab, 2 points to the 'Cross-validation' radio button, 3 points to the 'Classifier output' text area, and 4 points to the 'Result list' entry.

Interpretation of tree

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

Size of the tree : 9


Interpretation of tree

- ⦿ This indicates how the classifier uses the attributes to make a decision.
- ⦿ The leaf nodes indicate which class an instance will be assigned to should that node be reached.
- ⦿ The numbers in brackets after the leaf nodes indicate the number of instances assigned to that node, followed by how many of those instances are incorrectly classified as a result.

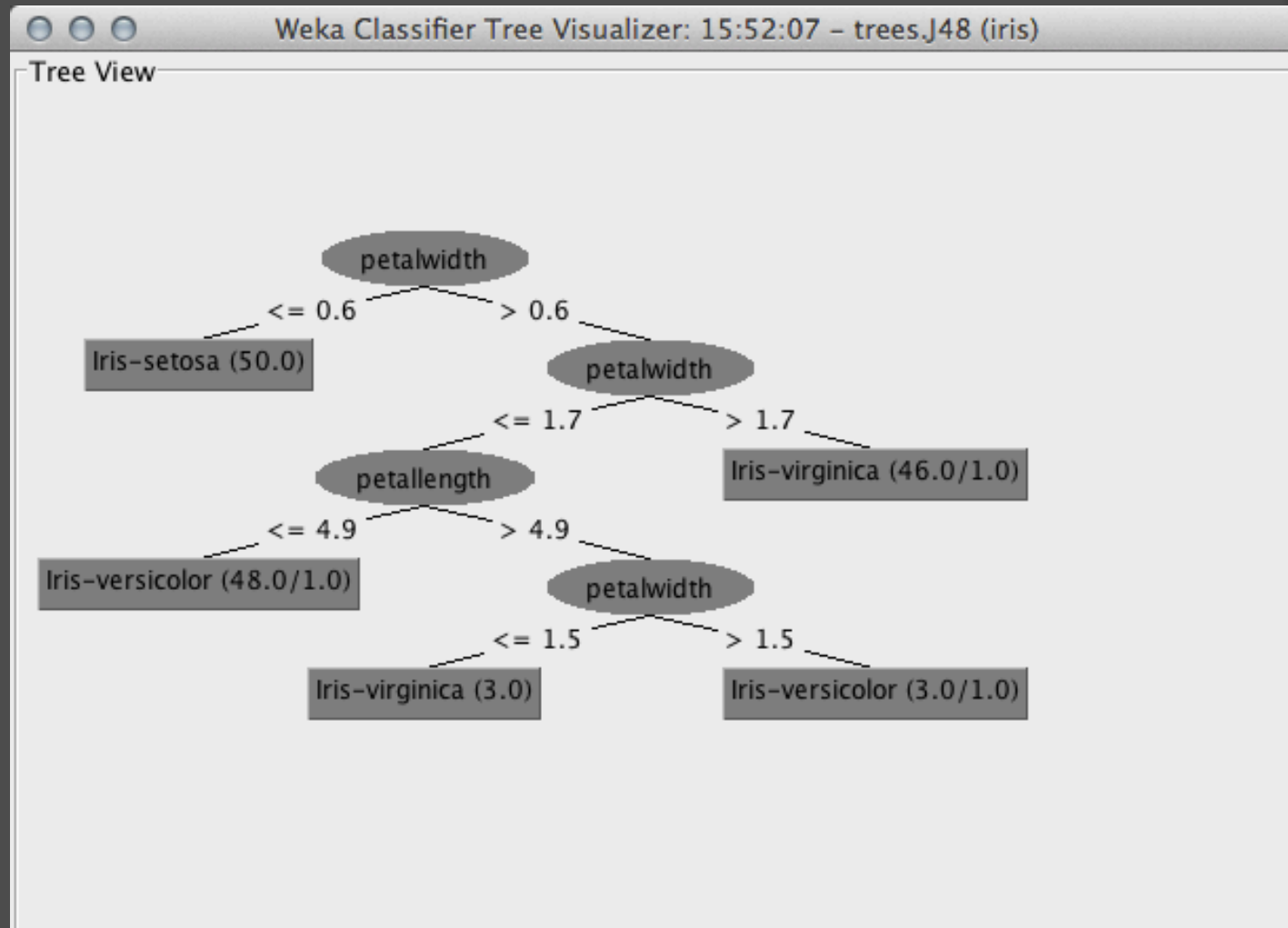
Tree visualization (1)

- View in main window
- View in separate window
- Save result buffer
- Delete result buffer

- Load model
- Save model
- Re-evaluate model on current test set
- Re-apply this model's configuration

- Visualize classifier errors
- Visualize tree 
- Visualize margin curve
- Visualize threshold curve ▶
- Cost/Benefit analysis ▶
- Visualize cost curve ▶

Tree visualization (2)



Exercise 1

- ⦿ Using the Weka explorer environment and load the training file “*diabetes.arff*” Perform classification with Naive Bayes, Decision Tree and K-NN (with $K=3$) Use the following setting :
 - 10 Fold Cross validation
 - 70% Training and 30% Test (percentage split)
- ⦿ Build a comparative table with the 2 different settings and the 3 classifiers and comment the results

Exercise 2

- Using the Weka explorer environment and load the training file “*diabetes.arff*” Perform classification with K-NN with different values of K (3,5,7,9,11,13) with 10 Fold Cross validation.
- Put the accuracy results in a table and comment the results. Emphasize how the results change in relation to the value of K

Exercise 3

- ⦿ Show the tree decision for “weather.arff” data using the following parameters:
 - Method: J48
 - Cross-validation: fixed on 5 and 10
- ⦿ Discuss the results (figure)

Clustering

Clustering

- ① The process of grouping physical or abstract objects into classes of similar objects i.e., given a set of records (instances, examples, objects, observations, ...), organize them into clusters (groups, classes)
- ① Works with both discrete and numerical data*

Classification vs clustering

- ***Classification: Supervised learning***

Learns a method for predicting the instance class from pre-labeled (classified) instances

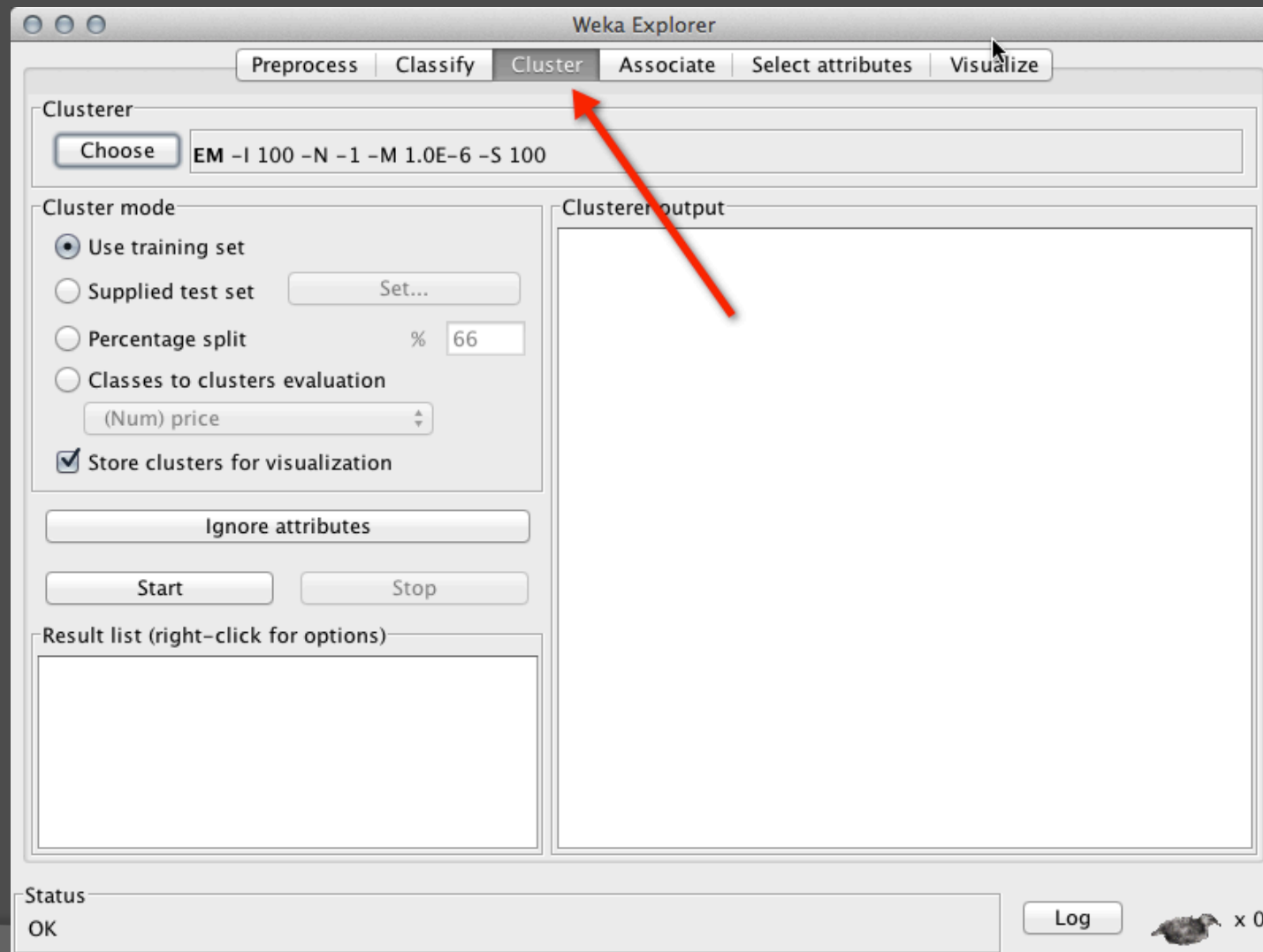
- ***Clustering: Unsupervised learning***

Finds “natural” grouping of instances given un-labeled data

Cluster definition

- ⦿ A cluster is a subset of objects which are “similar”
- ⦿ A subset of objects such that the distance between any two objects in the cluster is less than the distance between any object in the cluster and any object not located inside it
- ⦿ A connected region of a multidimensional space containing a relatively high density of objects

Clustering with Weka (1)



Clustering with Weka (2)

The screenshot shows the Weka Explorer application window with the 'Cluster' tab selected. The 'Clusterer' section is set to 'SimpleKMeans' with the following options: `-N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -S 10`. The 'Cluster mode' section has 'Use training set' selected, 'Percentage split' at 66%, and 'Store clusters for visualization' checked. The 'Clusterer output' area is empty. The 'Results list' area is also empty. The status bar at the bottom indicates 'Problem evaluating clusterer'.

Clusterer

Choose `SimpleKMeans -N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -S 10`

Cluster mode

- Use training set
- Supplied test set
- Percentage split %
- Classes to clusters evaluation
- Store clusters for visualization

Ignore attributes

Start Stop


Clusterer output

Results

Results list

Status

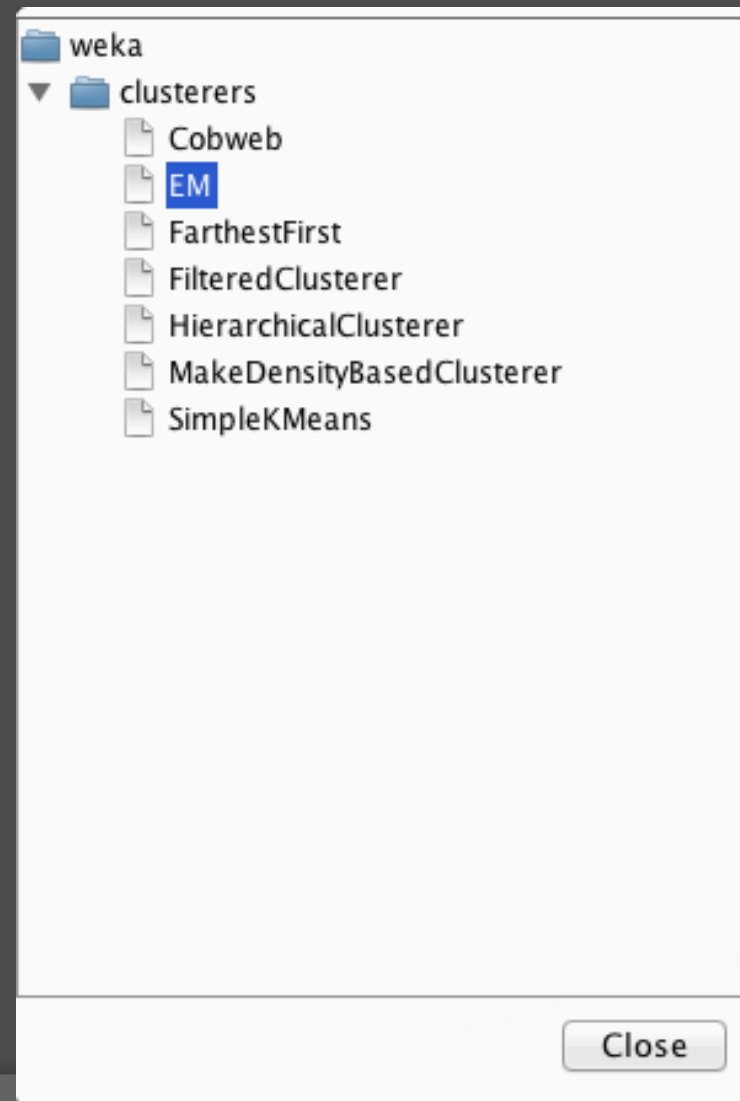
Problem evaluating clusterer

Log  x 1

Clustering: options

- ④ ***Use training set:*** After generating the clustering Weka classifies the training instances into clusters according to the cluster representation and computes the percentage of instances falling in each cluster
- ④ ***Supplied test set or Percentage split:*** Weka can evaluate clusterings on separate test data if the cluster representation is probabilistic (e.g. for EM).
- ④ ***Classes to clusters evaluation:*** In this mode Weka first ignores the class attribute and generates the clustering. Then during the test phase it assigns classes to the clusters, based on the majority value of the class attribute within each cluster (e.g. for k-Means)

Clustering with Weka (2)



EM

- ⦿ The EM clustering scheme generates probabilistic descriptions of the clusters in terms of mean and standard deviation for the **numeric attributes** and value **counts*** for the nominal ones

* incremented by 1 and modified with a small value to avoid zero probabilities

Example (1)

1. Open `auto_price.arff`
2. Go to *Clustering* tab
3. Choose a clustering method: *EM* 1
4. Set *Use training set* 2
5. Click on *Start* button 3

Example (2)

The screenshot shows the Weka Explorer interface with the 'Cluster' tab selected. The 'Clusterer' dropdown is set to 'EM -I 100 -N -1 -M 1.0E-6 -S 100', indicated by a red arrow and the number '1'. The 'Cluster mode' section has 'Use training set' selected, indicated by a red arrow and the number '2'. The 'Result list' shows '18:54:26 - EM' selected, indicated by a red arrow and the number '3'. The 'Clusterer output' pane displays the following information:

```
=== Run information ===  
Scheme:      weka.clusterers.EM -I 100 -N -1 -M 1.0E-6 -S 100  
Relation:    auto_price  
Instances:   159  
Attributes:  16  
             symboling  
             normalized-losses  
             wheel-base  
             length  
             width  
             height  
             curb-weight  
             engine-size  
             bore  
             stroke  
             compression-ratio  
             horsepower  
             peak-rpm  
             city-mpg  
             highway-mpg  
             price  
Test mode:   evaluate on training data  
  
=== Clustering model (full training set) ===  
  
EM  
==
```

The status bar at the bottom shows 'OK' and a 'Log' button.

Example (3)

Attribute	Cluster			
	0 (0.29)	1 (0.3)	2 (0.1)	3 (0.3)

symboling				
-3	1	1	1	1
-2	1	4	1	1
-1	1	13.0001	1.9999	8
0	10.2874	16.0629	5.7047	19.945
1	34.0247	4	2.9753	9
2	5	6.0038	10.0001	11.9961
3	1	10.9998	1	4.0002
[total]	53.3121	55.0666	23.6801	54.9413
normalized-losses				
mean	119.06	134.1044	100.9791	117.1393
std. dev.	28.2606	39.6996	20.3507	36.9184
wheel-base				
mean	94.3796	103.6076	94.3776	98.0115
std. dev.	0.9677	5.2916	3.9757	2.2365
length				
mean	161.729	185.184	161.6453	173.6787
std. dev.	5.2932	6.7599	10.7666	2.9438
width				
mean	63.9317	67.8757	64.2253	65.4333
std. dev.	0.2595	1.6323	1.3187	0.7726
height				
mean	52.9986	54.9413	53.9926	53.6924
std. dev.	2.1171	2.2688	2.0167	2.0217

Example (4)

```
mean                    5255.9762  5046.5452  4816.9411  5147.2922
std. dev.                355.5054   429.7359   401.8338   543.7446

city-mpg
mean                    30.9262    20.358    35.4016    25.3582
std. dev.                3.0629    3.0307    6.7859     1.614

highway-mpg
mean                    36.7524    25.518    41.1052    31.0112
std. dev.                3.1128    3.1878    7.46       2.2009

price
mean                   6959.1625 18703.0111  7412.919   9906.699
std. dev.              902.7482  5306.2492 1423.7368 1965.4932
```

Time taken to build model (full training data) : 1.88 seconds

=== Model and evaluation on training set ===

Clustered Instances

```
0      46 ( 29%)
1      48 ( 30%)
2      17 ( 11%)
3      48 ( 30%)
```



Log likelihood: -53.41921

Example (5)

The screenshot shows the Weka Explorer application window with the 'Cluster' tab selected. The 'Clusterer' section is set to 'EM -I 100 -N -1 -M 1.0E-6 -S 100'. The 'Cluster mode' section has 'Use training set' selected, and 'Store clusters for visualization' is checked. The 'Clusterer output' window displays the following text:

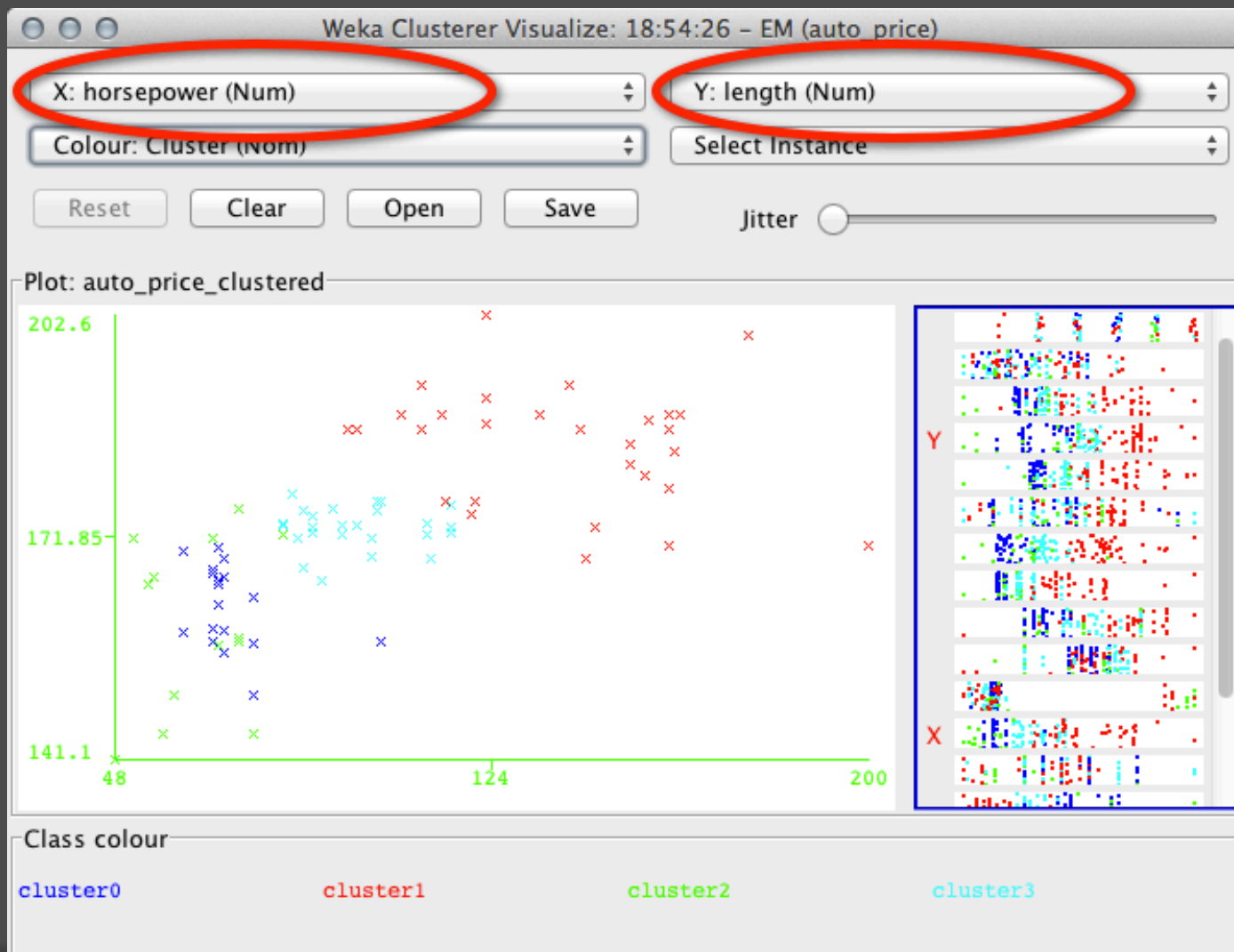
```
=== Run information ===  
Scheme:      weka.clusterers.EM -I 100 -N -1 -M 1.0E-6 -S  
Relation:    auto_price  
Instances:   159  
Attributes:  16  
             symboling  
             normalized-losses  
             wheel-base  
             length  
             width  
             height  
             curb-weight  
             engine-size  
             bore  
             stroke  
             compression-ratio  
             horsepower  
             peak-rpm  
             city-mpg  
             highway-mpg  
             price  
             evaluate on training data  
             training set) (full training set) ===
```

The 'Result list (right-click for options)' shows a single entry: '18:54:26 - EM'. A context menu is open over this entry, with the following options:

- View in main window
- View in separate window
- Save result buffer
- Delete result buffer
- Load model
- Save model
- Re-evaluate model on current test set
- Re-apply this model's configuration
- Visualize cluster assignments
- Visualize tree

A red arrow points to the 'Visualize cluster assignments' option. The status bar at the bottom left shows 'Status OK'.

Example (6)



Cobweb (1)

- Cobweb generates hierarchical clustering, where clusters are described probabilistically. The class attribute is ignored in order to allow later classes to clusters evaluation

Cobweb (2)

The screenshot shows the Weka Explorer application window with the 'Cluster' tab selected. The 'Clusterer' dropdown is set to 'Cobweb -A 1.0 -C 0.0028209479177387815 -S 42'. Under 'Cluster mode', 'Use training set' is selected, and 'Store clusters for visualization' is checked. The 'Ignore attributes' button is circled in red. The 'Clusterer output' pane shows the following information:

```
=== Run information ===
Scheme:      weka.clusterers.Cobweb -A 1.0 -C 0.0028209479
Relation:    weather
Instances:   14
Attributes:  5
              outlook
              temperature
              humidity
              windy
Ignored:
Test mode:   play
              evaluate on training data

=== Clustering model (full training set) ===
Number of merges: 1
Number of splits: 0
Number of clusters: 22

node 0 [14]
|
| node 1 [8]
| |
| | node 2 [2]
| | |
| | | leaf 3 [1]
| | | node 2 [2]
| | | |
| | | | leaf 4 [1]
| | node 1 [8]
| | |
| | | leaf 5 [1]
| | node 1 [8]
| | |
| | | leaf 6 [1]
```

The status bar at the bottom shows 'Status OK' and a 'Log' button.

Example (1)

1. Open weather.arff
2. Go to *Clustering* tab
3. Choose a clustering method: *Cobweb* 1
4. Set *Use training set* 2
5. Chose the “class” attribute on *Ignore Attribute* panel 3
6. Click on *Start* button 4

Example (2)

The screenshot shows the Weka Explorer interface with the 'Cluster' tab selected. The 'Clusterer' dropdown is set to 'Cobweb -A 1.0 -C 0.0028209479177387815 -S 42'. The 'Cluster mode' section has 'Use training set' selected, 'Percentage split' at 66%, and 'Store clusters for visualization' checked. The 'Clusterer output' pane shows run information and a clustering model. The 'Result list' shows '19:19:47 - Cobweb' selected. Red arrows and numbers 1-4 point to these specific elements.

1 Points to the Clusterer dropdown menu showing 'Cobweb -A 1.0 -C 0.0028209479177387815 -S 42'.

2 Points to the 'Use training set' radio button in the Cluster mode section.

3 Points to the 'Clustering model (full training set)' output in the Clusterer output pane.

4 Points to the '19:19:47 - Cobweb' entry in the Result list.

```
=== Run information ===
Scheme:      weka.clusterers.Cobweb -A 1.0 -C 0.0028209479177387815
Relation:    weather
Instances:   14
Attributes:  5
              outlook
              temperature
              humidity
              windy

Ignored:     play
Test mode:   evaluate on training data

=== Clustering model (full training set) ===
Number of merges: 1
Number of splits: 0
Number of clusters: 22

node 0 [14]
  node 1 [8]
    node 2 [2]
      leaf 3 [1]
    node 2 [2]
      leaf 4 [1]
  node 1 [8]
    leaf 5 [1]
  node 1 [8]
    leaf 6 [1]
```

Meaning of results (1)

```
=== Run information ===  
Scheme:          weka.clusterers.Cobweb -A 1.0 -C 0.00282094791773878:  
Relation:        weather  
Instances:       14  
Attributes:      5  
                 outlook  
                 temperature  
                 humidity  
                 windy  
Ignored:         play  
Test mode:       evaluate on training data  
  
=== Clustering model (full training set) ===  
  
Number of merges: 1  
Number of splits: 0  
Number of clusters: 22  
  
node 0 [14]  
| node 1 [8]  
| | node 2 [2]  
| | | leaf 3 [1]  
| | | node 2 [2]  
| | | | leaf 4 [1]  
| | node 1 [8]  
| | | leaf 5 [1]  
| | node 1 [8]  
| | | leaf 6 [1]
```

Meaning of results (2)

```
=== Clustering model (full training set) ===
```

```
Number of merges: 1  
Number of splits: 0  
Number of clusters: 22
```

```
node 0 [14]  
| node 1 [8]  
| | node 2 [2]  
| | | leaf 3 [1]  
| | | node 2 [2]  
| | | | leaf 4 [1]  
| | node 1 [8]  
| | | leaf 5 [1]  
| | node 1 [8]  
| | | leaf 6 [1]  
| | node 1 [8]  
| | | node 7 [3]  
| | | | leaf 8 [1]  
| | | node 7 [3]  
| | | | leaf 9 [1]  
| | | node 7 [3]  
| | | | leaf 10 [1]  
| | node 1 [8]  
| | | leaf 11 [1]  
node 0 [14]  
| node 12 [6]  
| | node 13 [2]  
| | | leaf 14 [1]  
| | node 13 [2]  
| | | leaf 15 [1]
```



Meaning of results (3)

- ⦿ Node N or leaf N represents a subcluster, whose parent cluster is N
- ⦿ The clustering tree structure is shown as a *horizontal tree*, where subclusters are aligned at the same column
- ⦿ The *root cluster is 0*. Each line with node 0 defines a subcluster of the root

Meaning of results (4)

- The number in *square brackets* after node N represents the number of instances in the parent cluster N
- Clusters with [1] at the end of the line are *instances*
- To view the clustering tree right click on the last line in the result list window and then select *Visualize tree*

Exercise 1

- ⦿ *Right click* on the last line in the result list window
- ⦿ Visualize cluster assignments - you get the Weka cluster *visualize* window
- ⦿ Put *Instance_number* on *X* and *Cluster* on *Y*
- ⦿ Click on *Save* and choose a file name (*.arff)
- ⦿ Explore the *arff* file and comment

k-Means

- ⦿ “ k ” stands for number of clusters, it is typically a user input to the algorithm; some criteria can be used to automatically estimate k
- ⦿ Works only for numerical data

Example (1)

1. Open weather.arff
2. Go to *Clustering* tab
3. Choose a clustering method:
SimpleKMeans 1
4. Set *Use training set* 2
5. Set *numCluster* (k) to 4 3
6. Click on *Start* button 4

Example (2)

The screenshot shows the Weka Explorer interface with the 'Cluster' tab selected. The 'Clusterer' dropdown is set to 'SimpleKMeans -N 4 -A "weka.core.EuclideanDistance -R first-last" -I 500 -S 10'. The 'Cluster mode' section has 'Use training set' selected. The 'Clusterer output' pane shows the following table:

Attribute	Full Data (768)	Cluster# 0 (173)	Cluster# 1 (222)	Cluster# 2 (36)	Cluster# 3 (337)
prex	3.8451	2.1214	7.7297	3.5556	2.2018
plas	120.8945	143.9191	130.0495	117	103.4599
pres	69.1055	73.104	77.4865	0.6667	68.8427
skin	20.5365	35.0289	16.8964	2	17.4748
insu	79.7995	194.6879	59.8063	0.6944	42.4421
mass	31.9926	36.9064	32.3122	25.7639	29.9249
pedi	0.4719	0.6211	0.465	0.3932	0.4082
age	33.2409	29.8613	46.8874	30.4444	26.2849

Annotations in the image:

- 1: Points to the 'Clusterer' dropdown menu.
- 2: Points to the 'Use training set' radio button.
- 3: Points to the 'Clusterer output' pane.
- 4: Points to the 'Result list' showing '20:25:14 - SimpleKMeans' selected.

Example (3)

```
Number of iterations: 31
Within cluster sum of squared errors: 95.23652346839076
Missing values globally replaced with mean/mode

Cluster centroids:

Attribute      Full Data      Cluster#
                (768)          0           1           2           3
                (173)        (222)       (36)       (337)
-----
preg           3.8451         2.1214      7.7297     3.5556     2.2018
plas          120.8945      143.9191    130.0495   117         103.4599
pres          69.1055       73.104      77.4865    0.6667     68.8427
skin          20.5365       35.0289    16.8964    2           17.4748
insu          79.7995      194.6879    59.8063    0.6944     42.4421
mass          31.9926       36.9064    32.3122    25.7639    29.9249
pedi           0.4719        0.6211     0.465      0.3932     0.4082
age           33.2409       29.8613    46.8874    30.4444    26.2849
```

Centroids

```
Time taken to build model (full training data) : 0.05 seconds
```

```
=== Model and evaluation on training set ===
```

```
Clustered Instances
```

```
0      173 ( 23%)
1      222 ( 29%)
2       36 (  5%)
3      337 ( 44%)
```

Meaning of results

- The first column gives you the overall population centroid. The second to fifth columns give you the centroids for cluster 0 to 4, respectively. Each row gives the centroid coordinate for the specific dimension.

Exercise 1

- ⦿ Go to the WEKA explorer environment and load the training file *iris.arff*
- ⦿ Cluster the iris dataset using the *k-Means* clustering algorithm with $k=5$. Watch the result given by WEKA (Cobweb).

Exercise 2

- Cluster the “*iris.arff*” dataset using the k-Means Clustering algorithm with $k=3$, $k=4$ and $k=5$, with the same ten different value of the seed parameter. Use the option: Classes to cluster evaluation to evaluate the accuracy and store the results on an excel file. Compute the mean of the three different k values for the k-Means.

Association Rules

Association rules mining

- ④ Method for discovering interesting relations between variables in large databases
- ④ For example, the rule {onion, potatoes} → {burger} would indicate that if a customer buys onions and potatoes together, he is likely to also buy hamburger meat

Classification vs Association Rules

⦿ Classification

- Focus on one target field
- Specify class in all cases
- Measures: Accuracy

⦿ Association Rules

- Many target fields
- Applicable in some cases
- Measures: Support, Confidence, Lift

Association rules

- ⦿ Association rule $R : \text{itemset1} \Rightarrow \text{itemset2}$
 - Itemset1, itemset2 are disjoint and Itemset2 is non-empty
 - meaning: if transaction includes Itemset1 then it also has Itemset2
- ⦿ Example
 - $A, B \Rightarrow E, C$

Association rules with Weka (1)

The screenshot shows the Weka Explorer application window. The 'Associate' tab is selected, and a red arrow points to the 'Generate...' button. The interface displays the following information:

Current relation
Relation: vote
Instances: 435
Attributes: 17
Sum of weights: 435

Selected attribute
Name: crime
Missing: 17 (4%)
Distinct: 2
Type: Nominal
Unique: 0 (0%)

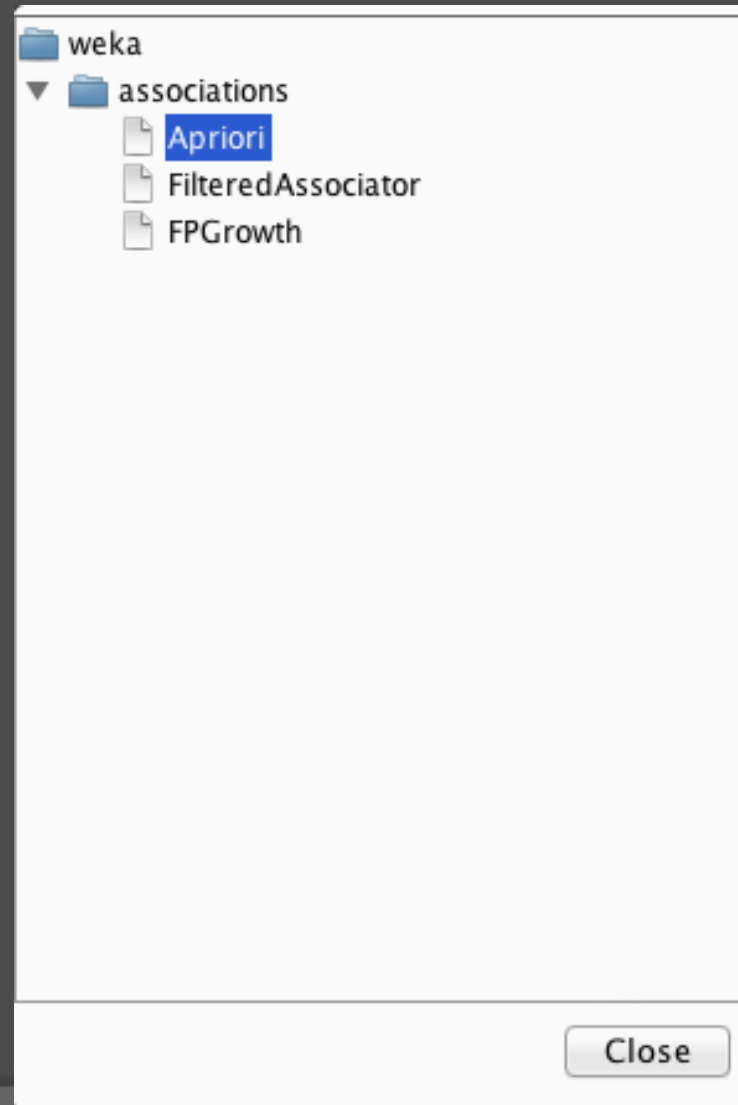
No.	Label	Count	Weight
1	n	170	170.0
2	y	248	248.0

Attributes
List of attributes: anti-satellite-test-ban, aid-to-nicaraguan-contras, mx-missile, immigration, synfuels-corporation-cutback, education-spending, superfund-right-to-sue, **crime**, duty-free-exports, export-administration-act-south-africa, Class.

Bar Chart
A bar chart showing the distribution of the 'crime' attribute. The 'n' category has a count of 170 (blue bar) and the 'y' category has a count of 248 (red bar).

Status
OK

Association rules with Weka (2)



Example (1)

1. Open `vote.arff`
2. Go to *Associate* tab
3. Choose a Association Rules method:
Apriori 1
4. Set *lowerMinBoundSupport* to 0.5 2
5. Set *numRules* to 15 3
6. Click on *Start* button 4

Example (2)

The screenshot shows the Weka Explorer interface with the 'Associate' tab selected. The 'Associator' section is set to 'Apriori' with the following command: `Apriori -N 15 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.5 -S -1.0 -c -1`. The 'Associator output' section displays the following text:

```
Apriori
=====
Minimum support: 0.45 (196 instances)
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=democr
2. adoption-of-the-budget-resolution=y physician-fee-freeze=n aid-to-nicaraguan-co
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)> lift:(1.62)
6. el-salvador-aid=n Class=democrat 200 ==> aid-to-nicaraguan-contras=y 197 <cc
7. el-salvador-aid=n 208 ==> aid-to-nicaraguan-contras=y 204 <conf:(0.98)> lift
8. adoption-of-the-budget-resolution=y aid-to-nicaraguan-contras=y Class=democrat
9. el-salvador-aid=n aid-to-nicaraguan-contras=y 204 ==> Class=democrat 197 <cc
10. aid-to-nicaraguan-contras=y Class=democrat 218 ==> physician-fee-freeze=n 210
```

Four red arrows point to specific elements: Arrow 1 points to the 'Choose' button; Arrow 2 points to the command line; Arrow 3 points to the 'Start' button; Arrow 4 points to the 'Result list' entry '18:31 14 - Apriori'.

Status: OK

Log  x 0

Meaning of results (1)

```
=== Run information ===  
  
Scheme:      weka.associations.Apriori -N 10 -T 1 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S  
Relation:    vote  
Instances:   435  
Attributes:  17  
             handicapped-infants  
             water-project-cost-sharing  
             adoption-of-the-budget-resolution  
             physician-fee-freeze  
             el-salvador-aid  
             religious-groups-in-schools  
             anti-satellite-test-ban  
             aid-to-nicaraguan-contras  
             mx-missile  
             immigration  
             synfuels-corporation-cutback  
             education-spending  
             superfund-right-to-sue  
             crime  
             duty-free-exports  
             export-administration-act-south-africa  
             Class  
  
=== Associator model (full training set) ===  
  
Apriori  
=====
```

Minimum support: 0.45 (196 instances)
Minimum metric <confidence>: 0.9

Meaning of results (2)

Apriori

Minimum support: 0.45 (196 instances)

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 210
2. adoption-of-the-budget-resolution=y physician-fee-freeze=n aid-to-nicaraguan-contras=y 211 ==> Class=democrat 210
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:(0.99)> lift:(1.62)
5. physician-fee-freeze=n 247 ==> Class=democrat 245 <conf:(0.99)> lift:(1.62)
6. el-salvador-aid=n Class=democrat 200 ==> aid-to-nicaraguan-contras=y 197 <conf:(0.98)> lift:(1.62)
7. el-salvador-aid=n 208 ==> aid-to-nicaraguan-contras=y 204 <conf:(0.98)> lift:(1.62)
8. adoption-of-the-budget-resolution=y aid-to-nicaraguan-contras=y Class=democrat 204
9. el-salvador-aid=n aid-to-nicaraguan-contras=y 204 ==> Class=democrat 197 <conf:(0.98)> lift:(1.62)
10. aid-to-nicaraguan-contras=y Class=democrat 218 ==> physician-fee-freeze=n 210

Meaning of results (3)

Apriori

Minimum support: 0.45 (196 instances)

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=democ
2. adoption-of-the-budget-resolution=y physician-fee-freeze=n aid-to-nicaraguan-c
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 <con
5. physician-fee-freeze=n 247 ==> Class=democrat 245 <conf:(0.99)> lift:(1.62)
6. el-salvador-aid=n Class=democrat 200 ==> aid-to-nicaraguan-contras=y 197 <c
7. el-salvador-aid=n 208 ==> aid-to-nicaraguan-contras=y 204 <conf:(0.98)> lif
8. adoption-of-the-budget-resolution=y aid-to-nicaraguan-contras=y Class=democrat
9. el-salvador-aid=n aid-to-nicaraguan-contras=y 204 ==> Class=democrat 197 <c
10. aid-to-nicaraguan-contras=y Class=democrat 218 ==> physician-fee-freeze=n 210

Exercise (1)

- ⦿ Mining the file 'supermarket.arff'
- ⦿ Open with a text editor this file and look at the value inside the file.
- ⦿ Which is the particularity of this file?
- ⦿ Try to understand why this file is particularly adapted for the Association Rules task.

Exercise (2)

- Create an “arff”-file containing the following document-word representation (binary mode).

t1 = {machine, learning, classifier}

t2 = {data, mining, associative, classifier}

t3 = {mining, decision, tree}

t4 = {association, mining, data}

t5 = {decision, tree, classifier}

Exercise (3)

- ⦿ Extract the top 10 Association Rules from your 'arff'-file (Exercise 2)
- ⦿ Discuss the results
- ⦿ Use a sparse mode representation and extract the association rules.

Questions concerning to
final project...