---

**Data Mining and Statistics Within the Health Services**

University of East Anglia

# Tutorial for Weka
## a data mining tool

Dr. Wenjia Wang
School of Computing Sciences
University of East Anglia

Data → Pre-processing → Data Mining → Knowledge

Data Mining & Statistics within the Health Services

---

# Content

University of East Anglia

1. Introduction to Weka
2. Data Mining Functions and Tools
3. Data Format
4. Hands-on Demos
   4.1 Weka Explorer
   - Classification
   - Attribute( feature) Selection
   4.2 Weka Experimenter
   4.3 Weka KnowledgeFlow
5. Summary

Data Mining & Statistics within the Health Services        Weka Tutorial (Dr. Wenjia Wang)        2

---

# 1. Introduction to WEKA

University of East Anglia

- A collection of open source of many data mining and machine learning algorithms, including
  - pre-processing on data
  - Classification:
  - clustering
  - association rule extraction
- Created by researchers at the University of Waikato in New Zealand
- Java based (also open source).

Data Mining & Statistics within the Health Services        Weka Tutorial (Dr. Wenjia Wang)        3

---

# Weka Main Features

University of East Anglia

- 49 data preprocessing tools
- 76 classification/regression algorithms
- 8 clustering algorithms
- 15 attribute/subset evaluators + 10 search algorithms for feature selection.
- 3 algorithms for finding association rules
- 3 graphical user interfaces
  - "The Explorer" (exploratory data analysis)
  - "The Experimenter" (experimental environment)
  - "The KnowledgeFlow" (new process model inspired interface)

Data Mining & Statistics within the Health Services        Weka Tutorial (Dr. Wenjia Wang)        4

---

# Weka: Download and Installation

University of East Anglia

- Download Weka (the stable version) from
  http://www.cs.waikato.ac.nz/ml/weka/
  - Choose a self-extracting executable (including Java VM)

  - (If you are interested in modifying/extending weka there is a developer version that includes the source code)

- After download is completed, run the self-extracting file to install Weka, and use the default set-ups.

Data Mining & Statistics within the Health Services        Weka Tutorial (Dr. Wenjia Wang)        5
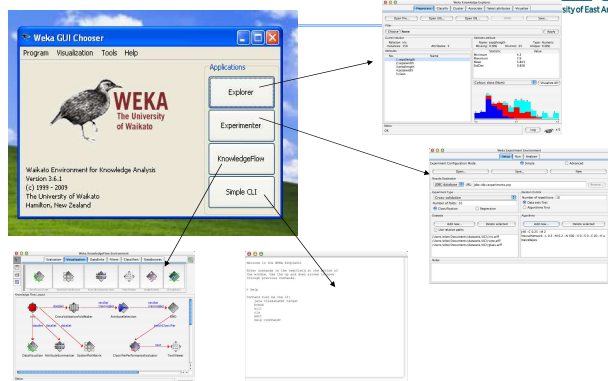
---

# Start the Weka

University of East Anglia

- From windows desktop,
  - click "Start", choose "All programs",
  - Choose "Weka 3.6" to start Weka
  - Then the first interface window appears:
    Weka **GUI Chooser**.



Data Mining & Statistics within the Health Services        Weka Tutorial (Dr. Wenjia Wang)        6

---

## WEKA Application Interfaces

## Weka Application Interfaces

- Explorer
  - preprocessing, attribute selection, learning, visualiation
- Experimenter
  - testing and evaluating machine learning algorithms
- Knowledge Flow
  - visual design of KDD process
  - Explorer
- Simple Command-line
  - A simple interface for typing commands

## 2. Weka Functions and Tools

- Preprocessing Filters
- Attribute selection
- Classification/Regression
- Clustering
- Association discovery
- Visualization

## Load data file and Preprocessing

- Load data file in formats: ARFF, CSV, C4.5, binary
- Import from URL or SQL database (using JDBC)
- Preprocessing filters
  - Adding/removing attributes
  - Attribute value substitution
  - Discretization
  - Time series filters (delta, shift)
  - Sampling, randomization
  - Missing value management
  - Normalization and other numeric transformations

## Feature Selection

- Very flexible: arbitrary combination of search and evaluation methods
- Search methods
  - best-first
  - genetic
  - ranking ...
- Evaluation measures
  - ReliefF
  - information gain
  - gain ratio
- Demo data: weather_nominal.arff

## Classification

- Predicted target must be categorical
- Implemented methods
  - decision trees(J48, etc.) and rules
  - Naïve Bayes
  - neural networks
  - instance-based classifiers …
- Evaluation methods
  - test data set
  - crossvalidation
- Demo data: iris, contact lenses, labor, soybeans, etc.

## Clustering

- Implemented methods
  - *k*-Means
  - EM
  - Cobweb
  - X-means
  - FarthestFirst…
- Clusters can be visualized and compared to "true" clusters (if given)
- Demo data:
  - any classification data may be used for clustering when its class attribute is filtered out.

Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          13

## Regression

- Predicted target is continuous
- Methods
  - linear regression
  - neural networks
  - regression trees …
- Demo data: cpu.arff,

Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          14

## Weka: Pros and cons

- pros
  - Open source,
    - Free
    - Extensible
    - Can be integrated into other java packages
  - GUIs (Graphic User Interfaces)
    - Relatively easier to use
  - Features
    - Run individual experiment, or
    - Build KDD phases
- Cons
  - Lack of proper and adequate documentations
  - Systems are updated constantly (Kitchen Sink Syndrome)

Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          15

## 3. WEKA data formats

- Data can be imported from a file in various formats:
  - ARFF (Attribute Relation File Format) has two sections:
    - the **Header** information defines attribute name, type and relations.
    - the **Data** section lists the data records.
  - CSV: Comma Separated Values (text file)
  - C4.5: A format used by a decision induction algorithm C4.5, requires two separated files
    - Name file: defines the names of the attributes
    - Date file: lists the records (samples)
  - binary
- Data can also be read from a URL or from an SQL database (using JDBC)

Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          16

## Attribute Relation File Format (arff)

An ARFF file consists of two distinct sections:

- the **Header** section defines attribute name, type and relations, start with a keyword.
  @**Relation** <data-name>
  @attribute <attribute-name> <type> or {range}
- the **Data** section lists the data records, starts with
  @**Data**
  list of data instances
- Any line start with % is the comments.

Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          17

## Breast Cancer data in ARFF

% Breast Cancer data*: 286 instances (no-recurrence-events: 201, recurrence-events: 85)
% Part 1: Definitions of attribute name, types and relations
@**relation breast-cancer**
  @attribute age {'10-19','20-29','30-39','40-49','50-59','60-69','70-79','80-89','90-99'}
  @attribute menopause {'lt40','ge40','premeno'}
  @attribute tumor-size {'0-4','5-9','10-14','15-19','20-24','25-29','30-34','35-39','40-44','45-49','50-54','55-59'}
  @attribute inv-nodes {'0-2','3-5','6-8','9-11','12-14','15-17','18-20','21-23','24-26','27-29','30-32','33-35','36-39'}
  @attribute node-caps {'yes','no'}
  @attribute deg-malig {'1','2','3'}
  @attribute breast {'left','right'}
  @attribute breast-quad {'left_up','left_low','right_up','right_low','central'}
  @attribute 'irradiat' {'yes','no'}
  @**attribute 'Class'** {'no-recurrence-events','recurrence-events'}

% Part 2: data section
@**data**
  '40-49','premeno','15-19','0-2','yes','3','right','left_up','no','recurrence-events'
  '50-59','ge40','15-19','0-2','no','1','right','central','no','no-recurrence-events'
  '50-59','ge40','35-39','0-2','no','2','left','left_low','no','recurrence-events'
  ……
* source: http://archive.ics.uci.edu/ml/datasets/Breast+Cancer

Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          18
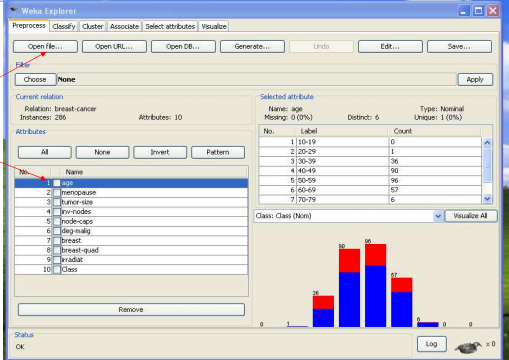
## 4.1 WEKA Explorer

UEA
University of East Anglia

- Click the Explorer on Weka GUI Chooser
- On the Explorer window,
  - click button "Open File" to open a data file from
    - the folder where your data files stored.
      e.g. Breast Cancer data: breast_cancer.arff
    Or (if you don't have this data set),
    - the data folder provided by the weka package:
      e.g. C:\Program Files\Weka-3-6\data
        using "iris.arff" or "weather_nominal.arff"

Data Mining & Statistics within the Health Services       Weka Tutorial (Dr. Wenjia Wang)       19

## Weka Explorer: open data file

UEA
University of East Anglia

- Open Breast Cancer data
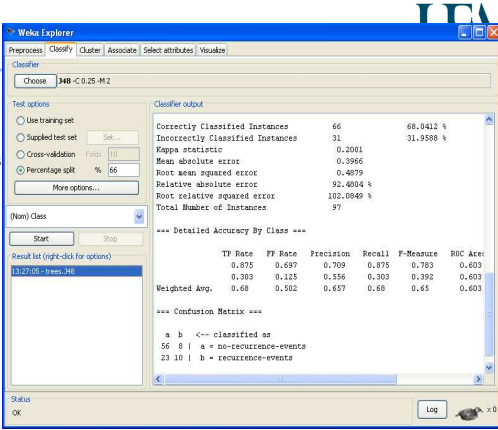- Click an attribute, e.g. age, then its distribution will be displayed in a histogram.



Data Mining & Statistics within the Health Services       Weka Tutorial (Dr. Wenjia Wang)       20

## Weka Explorer: training classifiers

UEA
University of East Anglia

After loaded a data file, click "Classify"

- Choose a classifier,
  - Under "Classifier": click "choose", then a drop-down menu appears,
  - Click "trees" and select "J48" – a decision tree algorithm
- Select a test option
  - Select "percentage split"
    - with default ratio 66% for training and 34% for testing
- Click "Start" to train and test the classifier.
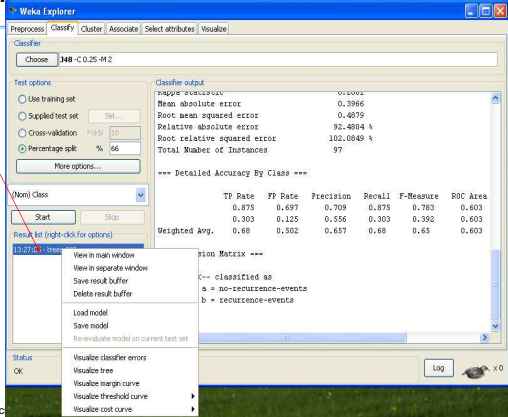  - The training and testing information will be displayed in classifier output window.

Data Mining & Statistics within the Health Services       Weka Tutorial (Dr. Wenjia Wang)       21

## Results

- Testing results:
- 97 cases used in test.

Correct:
66 (68%)

Wrong:
31 (32%)



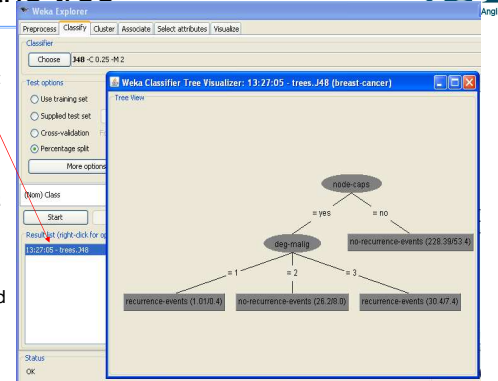Data Mining & Statistics within the Health Services       Weka Tutorial (Dr. Wenjia Wang)       22

## Options for results and model

UEA

- Point to result list window, and right click mouse.
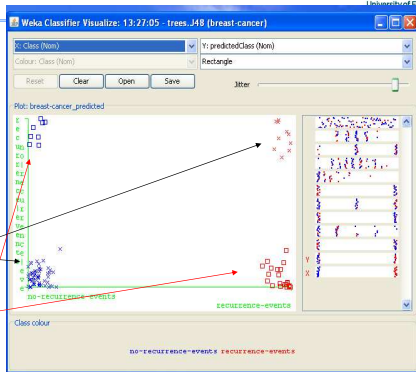- A menu will pop out to show all the options available about the model.



Data Mining & Statistics

## View the tree

UEA
Anglia

- Point to result list window, and right click mouse,
- Choose "**visualize tree**", then the tree will be displayed in another window.



Data Mining & Statistics within the Health Services       Weka Tutorial (Dr. Wenjia Wang)       24
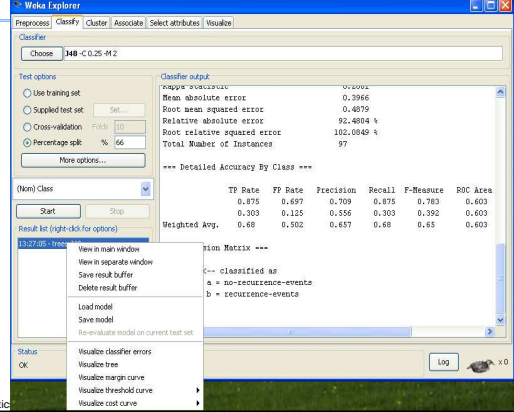
## View classifier errors

- right click the result list,
- Choose "**visualize classifier error**", then a new window will be popped out to display the classifier's error.
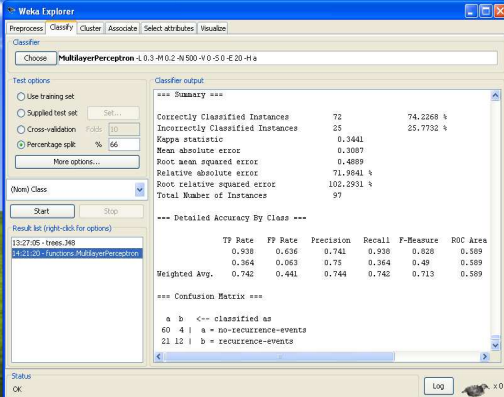  - Correctly predicted cases
  - Wrong cases



Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          25

## Save the model and results

- Right click on the result list
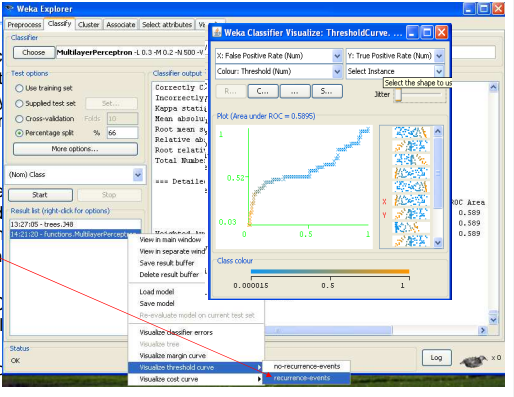- Choose "save model" and "save result buffer" to save the classifier and the results to the disk folder.



Data Mining & Statistic

## Train a neural net

Click "Choose" to select another function, e.g. "Multilayer Perceptron" - a type of neural net.

Then click "Start" to train and test it. (note: the training may take much longer time.)
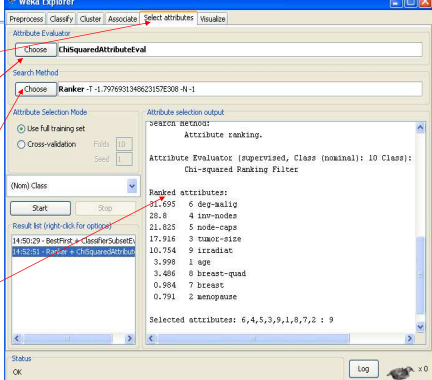
The results seem better than the tree classifier.



Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          27

## View the model's ROC curve

- Right click the result "Multilayer Perceptron"
- Choose "visualize threshold curve" and "recurrence-events";
- The ROC curve will be displayed



Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          28
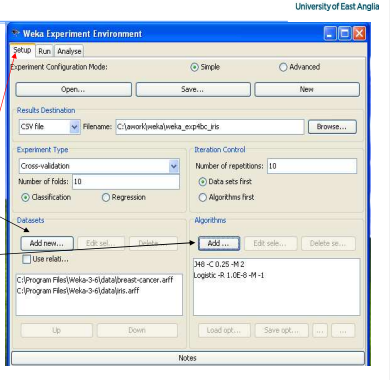
## Select Attributes

- Click "Select Attributes"
- Choose an "attribute evaluator"
  - e.g. chiSquare
- Choose a "Search Method"
- Then click "Start"
- The selected attributes are listed.



Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          29
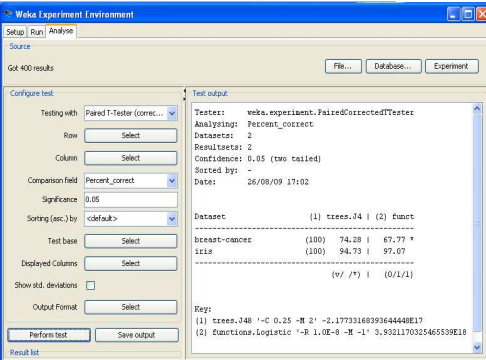
## 4.2 Weka Experimenter

- you can use Experimenter to carry out experiments for multiple data sets using multiple methods, e.g. classifying
- two data sets
  - Breast cancer
  - Iris
- Using two methods
  - Decision Tree: J48
  - Logistic
- The experiment is "Setup" as shown in the screenshot.
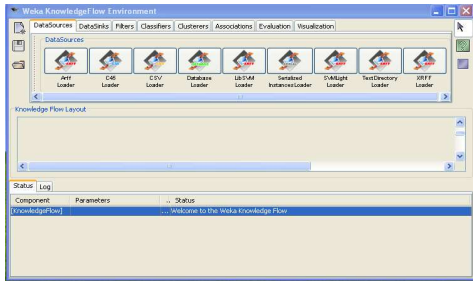- Then click "Run"



Data Mining & Statistics within the Health Services          Weka Tutorial (Dr. Wenjia Wang)          30

## Analysis of the results

- Click "analysis" to analyse the results,

  E.g.

  paired t-test significance
- Click "Experiment"
- Configure test: choosing appropriate test and parameters
- Click "Perform test" and the test results are listed.

## 4.3 KnowledgeFlow

- Click KnowledgeFlow on Weka GUI Chooser
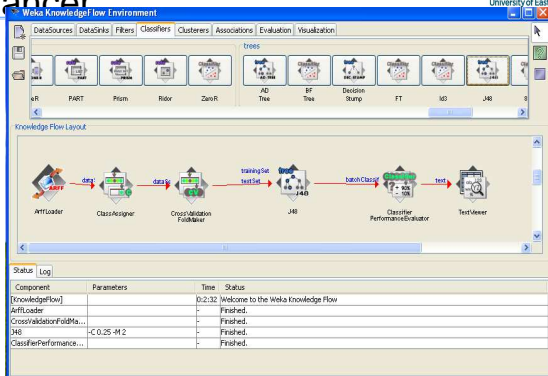- A new window opened for buidling KDD process.

## Steps for building a KDD process
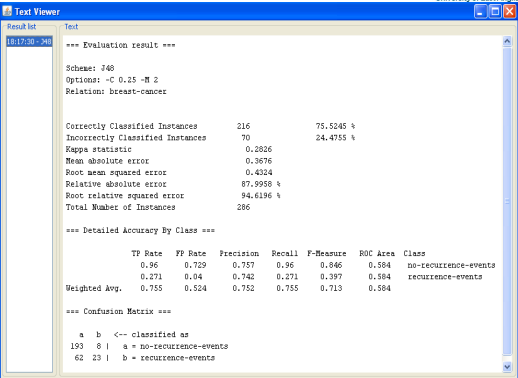
Major steps for building a process
1. Adding required nodes
   1) Add nodes
   2) Add a data source node from "DataSources"
      1) Right click to configure it with a data set
   3) Add a classAssigner node from "Evaluation" and a CrossValidationFoldmaker node
   4) Add a classifier, e.g. J48, from Classifiers
   5) Add a classiferPerformanceEvaluator node from "Evaluation"
   6) Add a text viewer from "Visualisation"
2. Connect the nodes
   – Right click "DataSource" node and choose DataSet, then connect it to the ClassAssigner node,
   – do the same or similar for connecting between the other nodes.
3. Run the process (using the default setups for each node)
   – Right click DataSource node and choose "Start loading", the process should run and "Status" window should indicate if the run is correct and completed.
4. View the results:
   – If the run is correctly completed, right click "Text Viewer" node and choose "Show results", then another window pops out to show the results.

## A KDD process for Breast Cancer

## Results of the KDD process

- right click "Text Viewer" node and choose "Show results" then another window pops out to show the results.

## 5. Weka Tutorial Summary

Weka is open source data mining software that offers
- Some GUI interfaces for data mining
  – Explorer
  – Experimenter
  – KnowledgeFlow
- Many functions and tools that include
  – Methods for **classification:**
     decision trees, rule learners, naive Bayes, decision tables, locally weighted regression, SVMs, instance-based learners, logistic regression, multi-layer perceptron
  – methods for **regression/prediction:**
     linear regression, model tree generators, locally weighted regression, instance-based learners, decision tables, multi-layer perceptron
  – **Ensemble schemes**
     • Bagging, boosting, stacking, RandomFrest
  – Methods for **clustering**:
     • K-means, EM and Cobweb
  – Methods for feature selection