# Fondamenti di Data Science e Machine Learning
## WEKA

# Outline

▸ **Introduction to WEKA**

▸ **The Explorer**

  ▸ Process data

  ▸ Classification

  ▸ Clustering

  ▸ Association Rules

  ▸ Attribute Selection

  ▸ Data Visualization

▸ **References and Resources**

# WEKA

▸ Waikato Environment for Knowledge Analysis

> ▸ It's a data mining/machine learning tool developed by Department of Computer Science, University of Waikato, New Zealand

> ▸ Weka is also a bird found only on the islands of New Zealand

# How to install WEKA

▸ Download and Install WEKA

  ▸ Website:
  ▸ https://waikato.github.io/weka-wiki/downloading_weka/
  ▸ Support multiple platforms (written in java):
    ▸ Windows, Mac OS X and Linux

| Project | Software | Book | Courses | Publications | People | Related |
|---|---|---|---|---|---|---|

## Downloading and installing Weka

There are two versions of Weka: Weka 3.8 is the latest stable version and Weka 3.9 is the development version. For the bleeding edge, it is also possible to download nightly snapshots. Stable versions receive only bug fixes, while the development version receives new features.

Weka 3.8 and 3.9 feature a package management system that makes it easy for the Weka community to add new functionality to Weka. The package management system requires an internet connection in order to download and install packages.

# WEKA: Features

▸ Main Features

- ▸ 49 data preprocessing tools

- ▸ 76 classification/regression algorithms

- ▸ 8 clustering algorithms

- ▸ 3 algorithms for finding association rules

- ▸ 15 attribute/subset evaluators + 10 search algorithms for feature selection

# The Graphical User Interface of WEKA

▸ Three graphical user interfaces

  ▸ "The Explorer" (exploratory data analysis)

  ▸ "The Experimenter" (experimental environment)

  ▸ "The KnowledgeFlow" (new process model inspired interface)

# The Explorer

▸ Data can be imported from a file in various formats: ARFF, CSV, C4.5, binary

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

▸ Pre-processing tools in WEKA are called "filters"

▸ WEKA contains filters for:

  ▸ Discretization, normalization, resampling, attribute selection, transforming and combining attributes, …

# Data Preprocessing (1)

▶ WEKA only deals with "flat" files

@relation cancerTraining

@attribute   'ClumpThickness'         NUMERIC

@attribute   'UniformityofCellSize'   NUMERIC

@attribute   'UniformityofCellShape'  NUMERIC

@attribute   'MarginalAdhesion'       NUMERIC

@attribute   'SingleEpithelialCellSize'  NUMERIC

@attribute   'BareNuclei' {1,2,3,4,5,6,7,8,9,10,'?'}

@attribute   'BlandChromatin'         NUMERIC

@attribute   'NormalNucleoli'         NUMERIC

@attribute   'Mitoses'     NUMERIC

@attribute   'Class'         {2,4}

Flat file in ARFF format

@data

5,1,1,1,2,1,3,1,1,2

5,4,4,5,7,10,3,2,1,4

….

# Data Preprocessing (1)

▸ WEKA only deals with "flat" files

@relation cancerTraining

@attribute  'ClumpThickness'         NUMERIC         ⟵  **Numeric attribute**

@attribute  'UniformityofCellSize'   NUMERIC

@attribute  'UniformityofCellShape'  NUMERIC

@attribute  'MarginalAdhesion'       NUMERIC

@attribute  'SingleEpithelialCellSize'  NUMERIC

@attribute  'BareNuclei' {1,2,3,4,5,6,7,8,9,10,'?'}  ⟵  **Nominal attribute**

@attribute  'BlandChromatin'         NUMERIC

@attribute  'NormalNucleoli'         NUMERIC

@attribute  'Mitoses'     NUMERIC

@attribute  'Class'          {2,4}

@data

5,1,1,1,2,1,3,1,1,2

5,4,4,5,7,10,3,2,1,4

….

# Data Preprocessing (2)

# Data Preprocessing (2)

# Data Preprocessing (3)

# Data Preprocessing (3)

# Data Preprocessing (4)

# Data Preprocessing (4)

# Data Preprocessing (5)

# Classification

▶ Explorer: building "classifiers"

  ▶ Classifiers in WEKA are models for predicting nominal or numeric quantities

  ▶ Implemented learning schemes include:

    ▸ Decision trees and lists

    ▸ instance-based classifiers

    ▸ support vector machines

    ▸ multi-layer perceptrons

    ▸ logistic regression

    ▸ Bayes' nets

    ▸ etc

# Decision Tree

▸ **Algorithm for Decision Tree Induction**

▸ **Basic algorithm (a greedy algorithm)**

- ▸ The tree is constructed in a top-down recursive divide-and-conquer
- ▸ At the start, all the training examples are at the root
- ▸ Attributes are categorical (if continuous-valued, they are discretized in advance)
- ▸ Examples are partitioned recursively based on selected attributes
- ▸ Test attributes are selected on the basis of a heuristic or statistical measure (e.g., information gain)

# Classification: An example (1)

# Classification: An example (1)

# Classification: An example (1)

# Classification: An example (1)

# Classification: An example (1)

# Classification: An example (2)

# Classification: An example (2)

# Classification: An example (3)

# Classification: An example (3)

# Classification: An example (3)

# Classification: An example (4)

# Classification: An example (4)

# Classification: An example (4)

# Classification: An example (4)

# Classification: An example (5)

# Classification: An example (5)

# Classification: An example (5)

# Classification: An example (5)

# Classification: An example (5)

# Classification: An example (6)

# Classification: An example (6)

# Classification: An example (7)

# Clustering (1)

▸ WEKA contains "clusters" for finding groups of similar instances in a dataset

▸ Implemented schemes are:

  ▸ k-Means, EM, Cobweb, X-means, FarthestFirst

▸ Clusters can be visualized and compared to "true" clusters (if given)

▸ Evaluation based on loglikelihood if clustering scheme produces a probability distribution

# Clustering (2)

▸ **The K-Means Clustering Method**

▸ Given k, the k-means algorithm is implemented in four steps:

1. Partition objects into k nonempty subsets

2. Compute seed points as the centroids of the clusters of the current partition (the centroid is the centre, i.e., mean point, of the cluster)

3. Assign each object to the cluster with the nearest seed point

4. Go back to Step 2, stop when no more new assignment

# Clustering: An example (1)

# Clustering: An example (1)

# Clustering: An example (2)

# Clustering: An example (2)

# of chosen clusters

# Clustering: An example (4)

# Clustering: An example (5)

# Clustering: An example (5)

# Clustering: An example (5)

# Clustering: An example (6)

# Clustering: An example (6)

# Clustering: An example (7)

# Associations Rules (1)

▸ WEKA contains an implementation of the Apriori algorithm for learning association rules

  ▸ Works only with discrete data

▸ Can identify statistical dependencies between groups of attributes:

  ▸ Es. milk, butter $\Rightarrow$ bread, eggs (with confidence 0.9 and support 2000)

▸ Apriori can compute all rules that have given minimum support and exceed given confidence

# Associations Rules: An example (1)

# Associations Rules: An example (1)

# Associations Rules: An example (1)

# Associations Rules: An example (2)

# Associations Rules: An example (2)

# Attribute selection

▸ Panel that can be used to investigate which (subsets of) attributes are the most predictive ones

▸ Attribute selection methods contain two parts:

  ▸ A search method: best-first, forward selection, random, exhaustive, genetic algorithm, ranking

  ▸ An evaluation method: correlation-based, wrapper, information gain, chi-squared, …

▸ Very flexible: WEKA allows (almost) arbitrary combinations of these two

# Attribute selection: An example (1)

# Attribute selection: An example (1)

# Attribute selection: An example (2)

# Attribute selection: An example (2)

# Attribute selection: An example (2)

# Attribute selection: An example (3)

# Attribute selection: An example (4)

# Data Visualization (1)

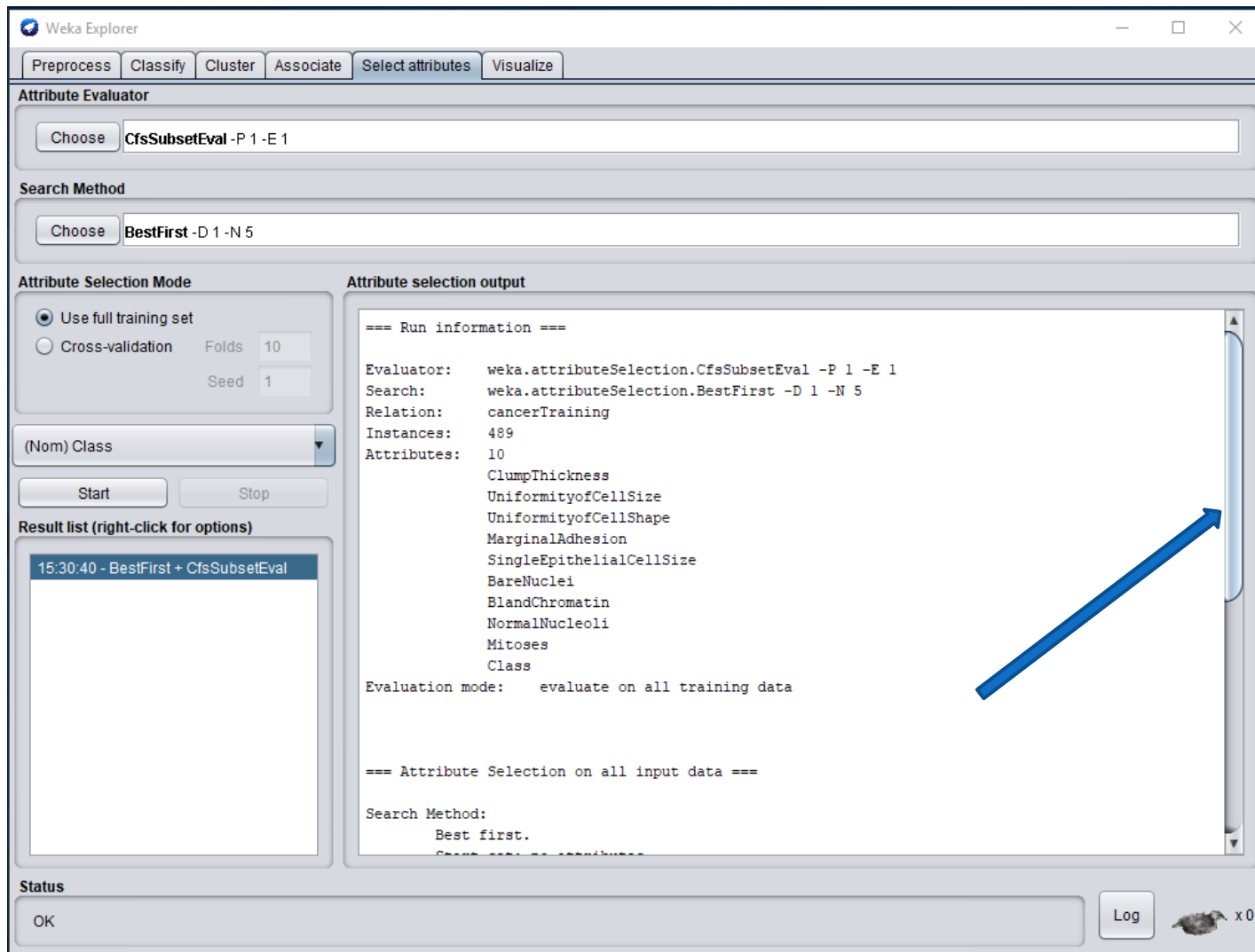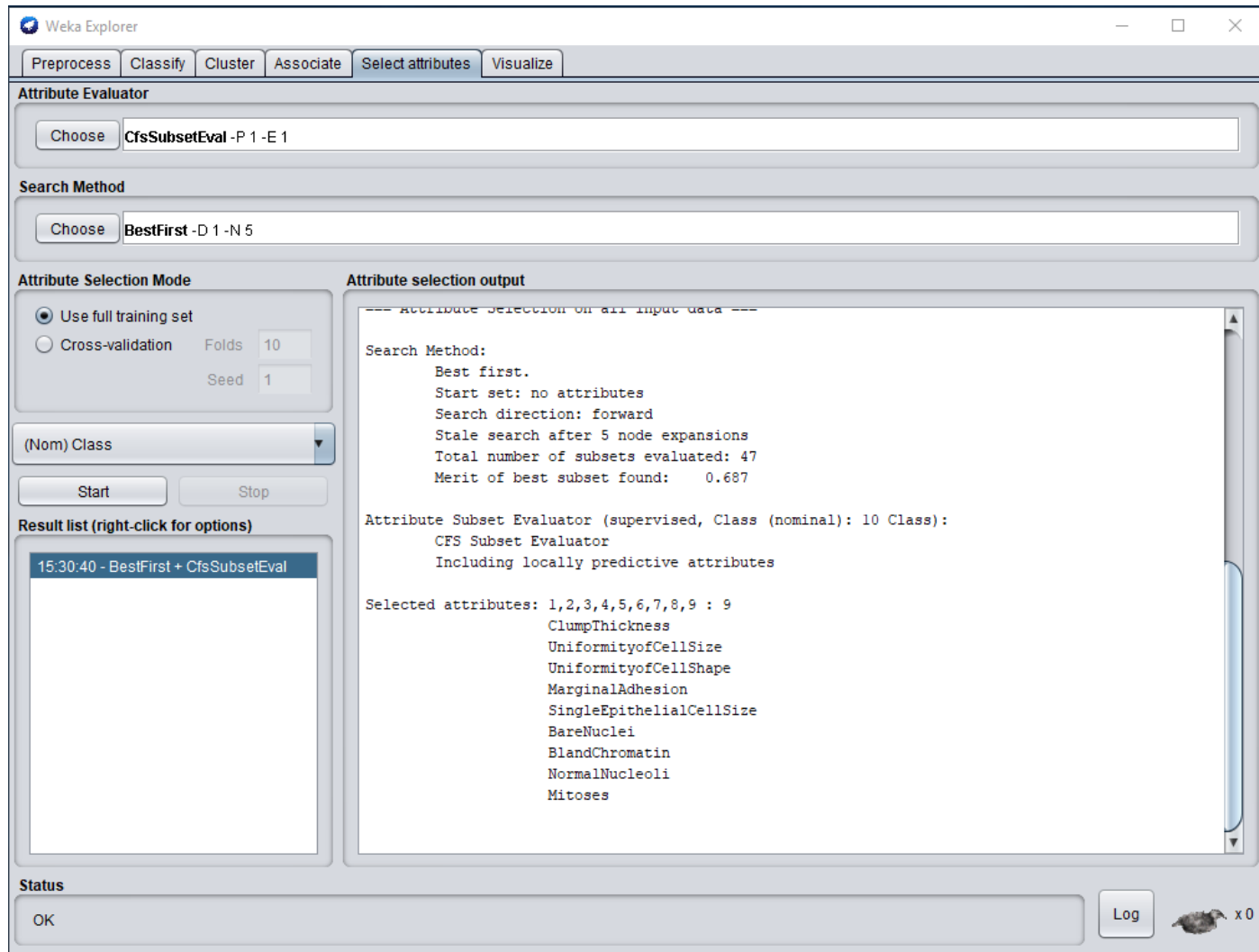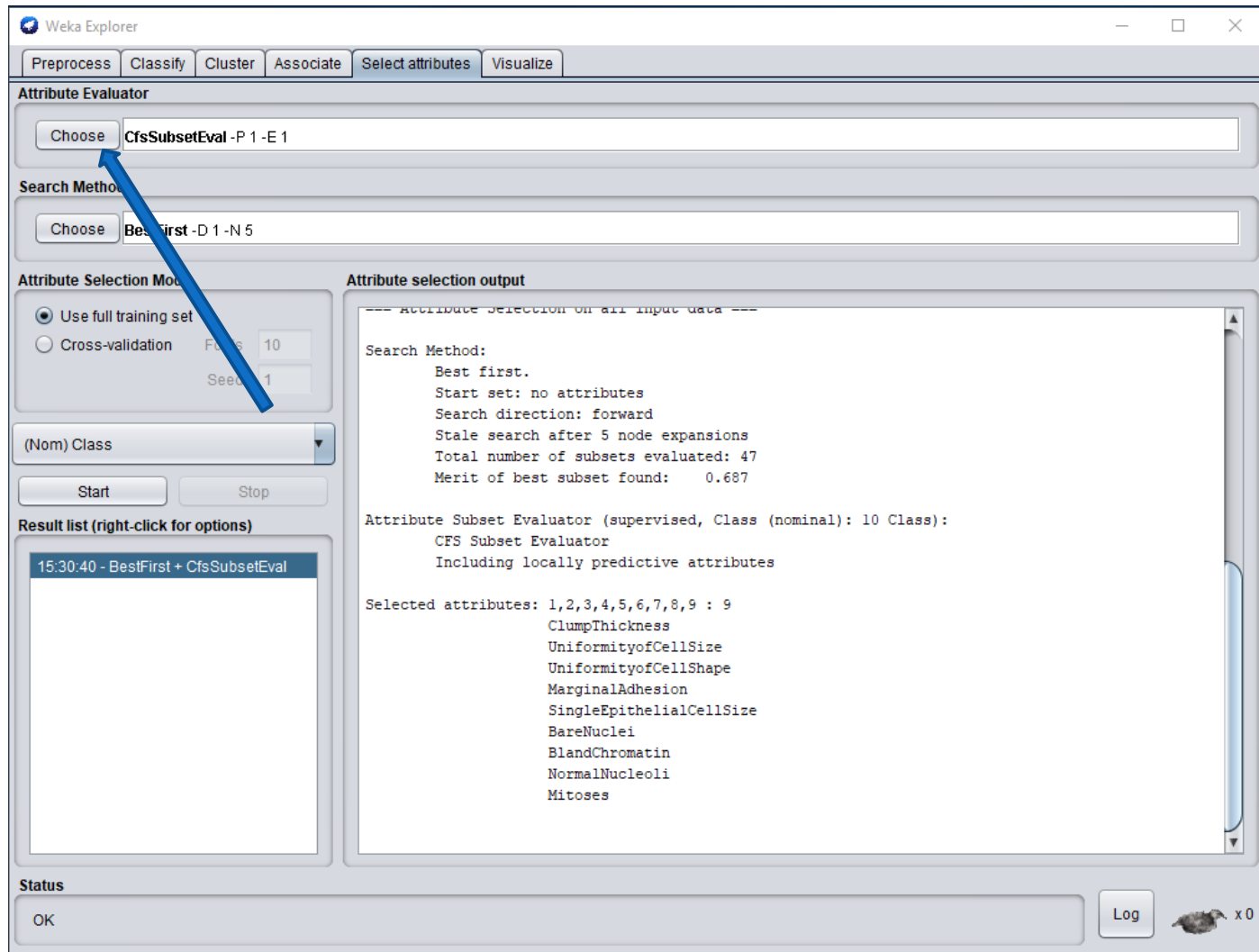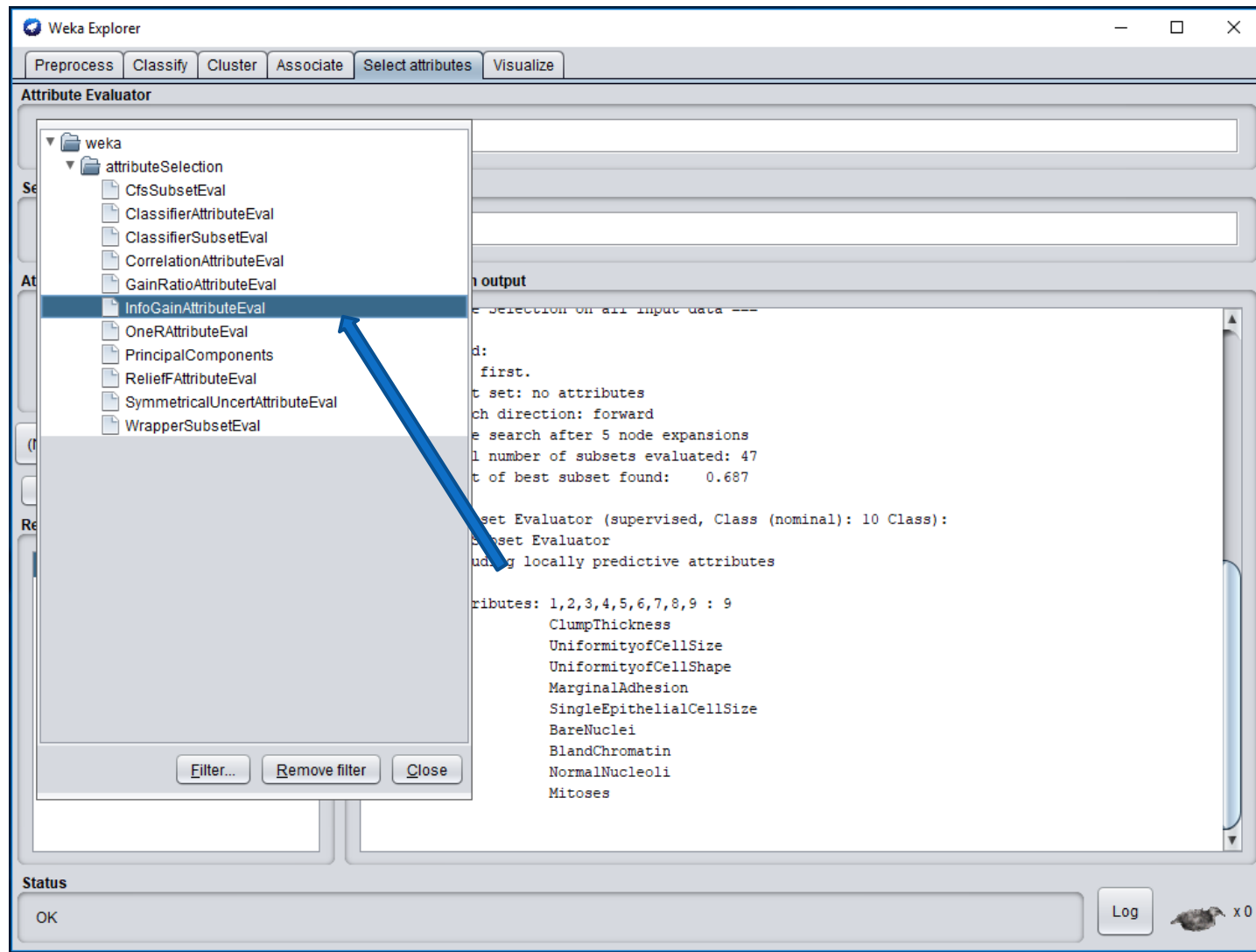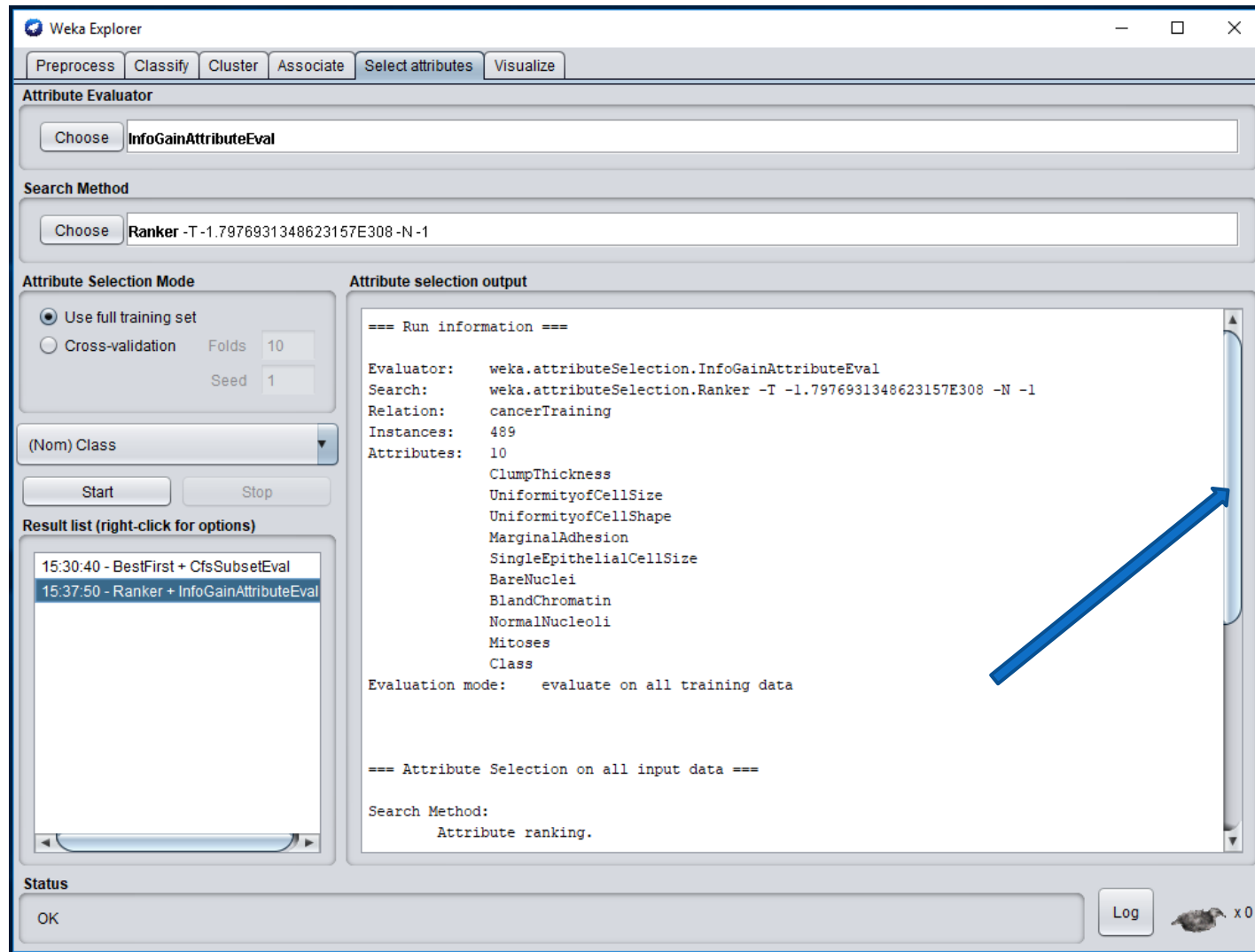▸ Visualization very useful in practice: e.g. helps to determine the difficulty of the learning problem

▸ WEKA can visualize single attributes (1-d) and pairs of attributes (2-d)

  ▸ To do: rotating 3-d visualizations (Xgobi-style)

▸ Colour-coded class values

▸ "Jitter" option to deal with nominal attributes (and to detect "hidden" data points)

▸ "Zoom-in" function

# Data Visualization (2)

# Data Visualization (2)

# Data Visualization (3)

# Data Visualization (3)

# Data Visualization (4)

# References and Resources

- References:

  - WEKA website:

    - http://www.cs.waikato.ac.nz/~ml/weka/index.html

  - WEKA Tutorial:

    - Machine Learning with WEKA: A presentation demonstrating all graphical user interfaces (GUI) in Weka A presentation which explains how to use Weka for exploratory data mining
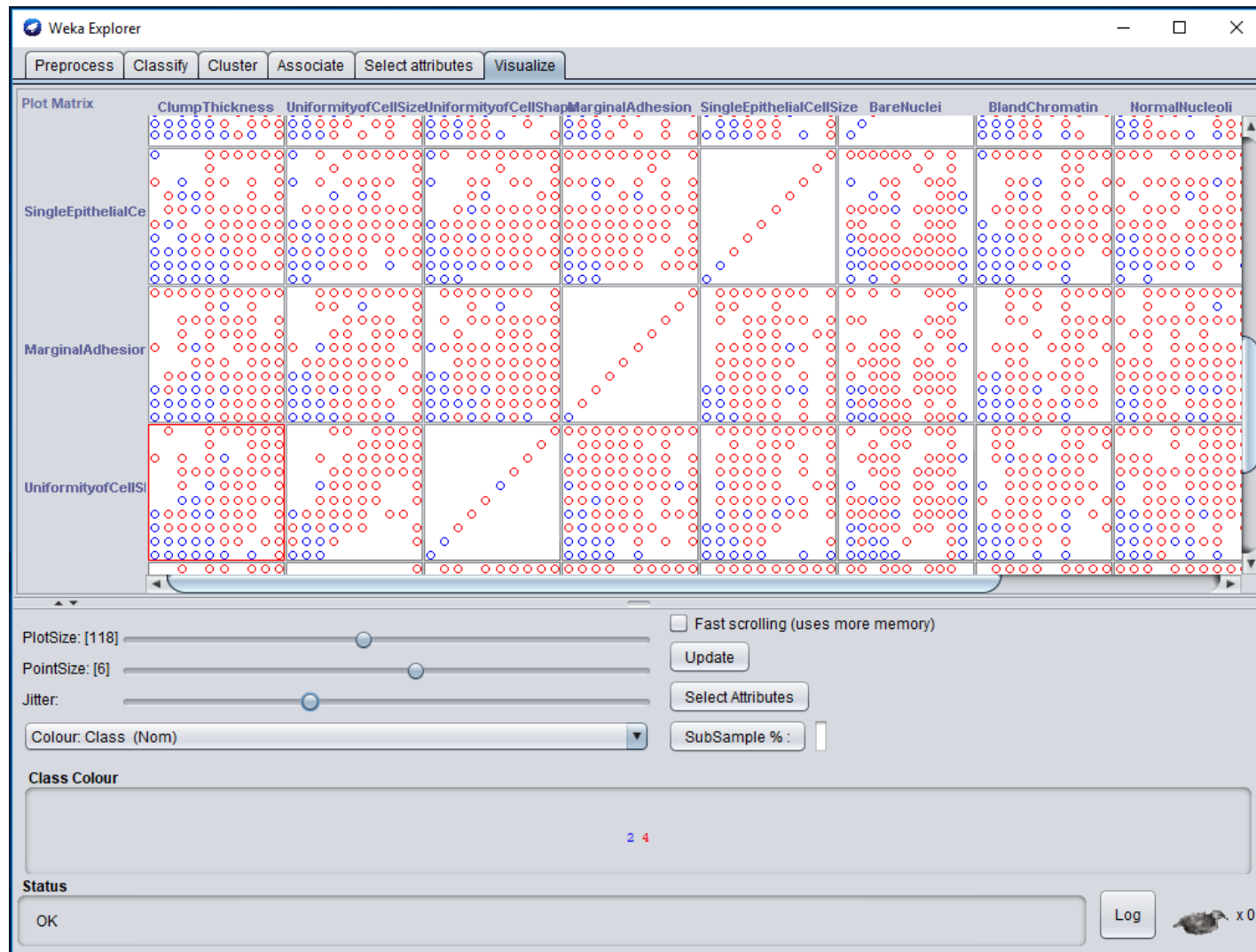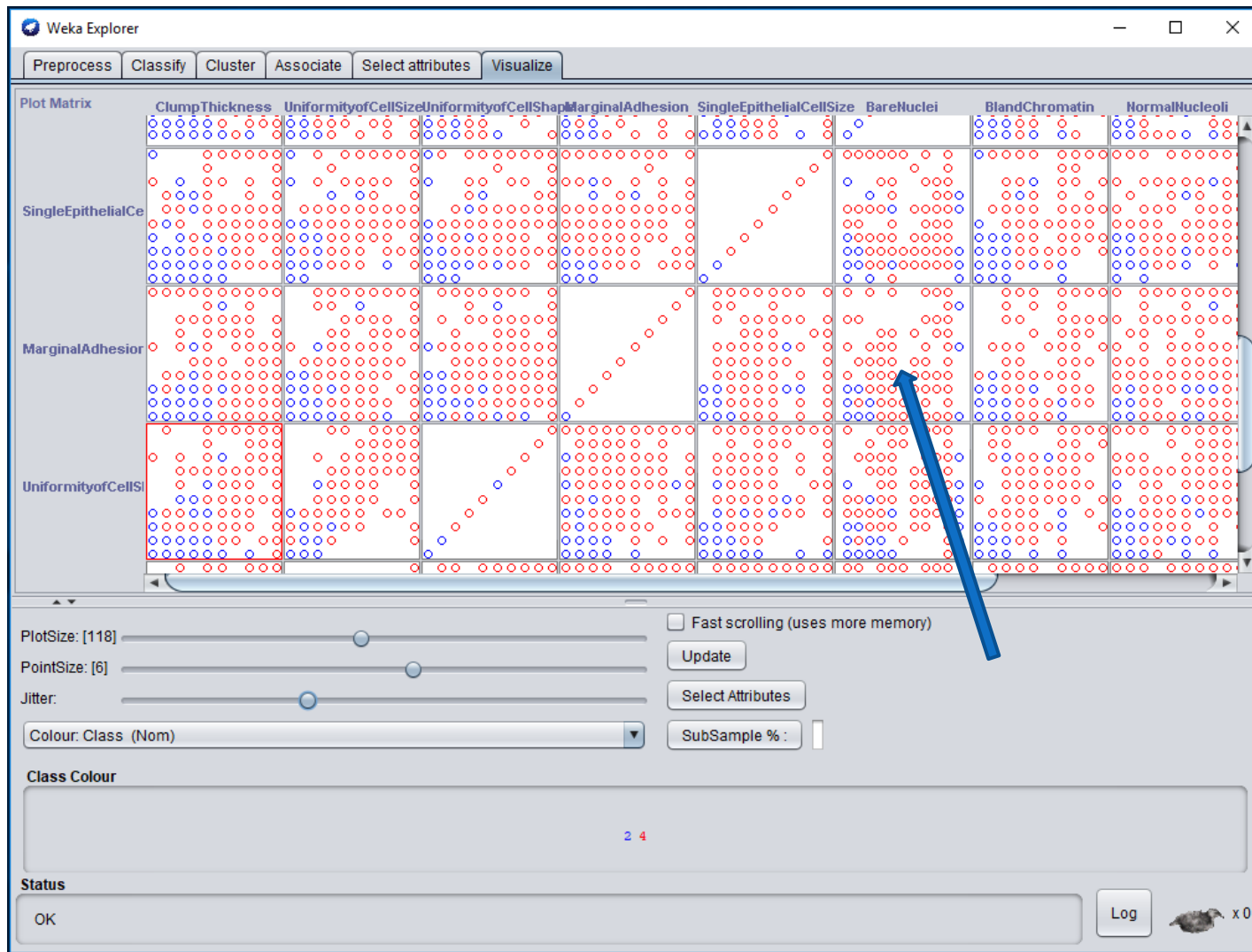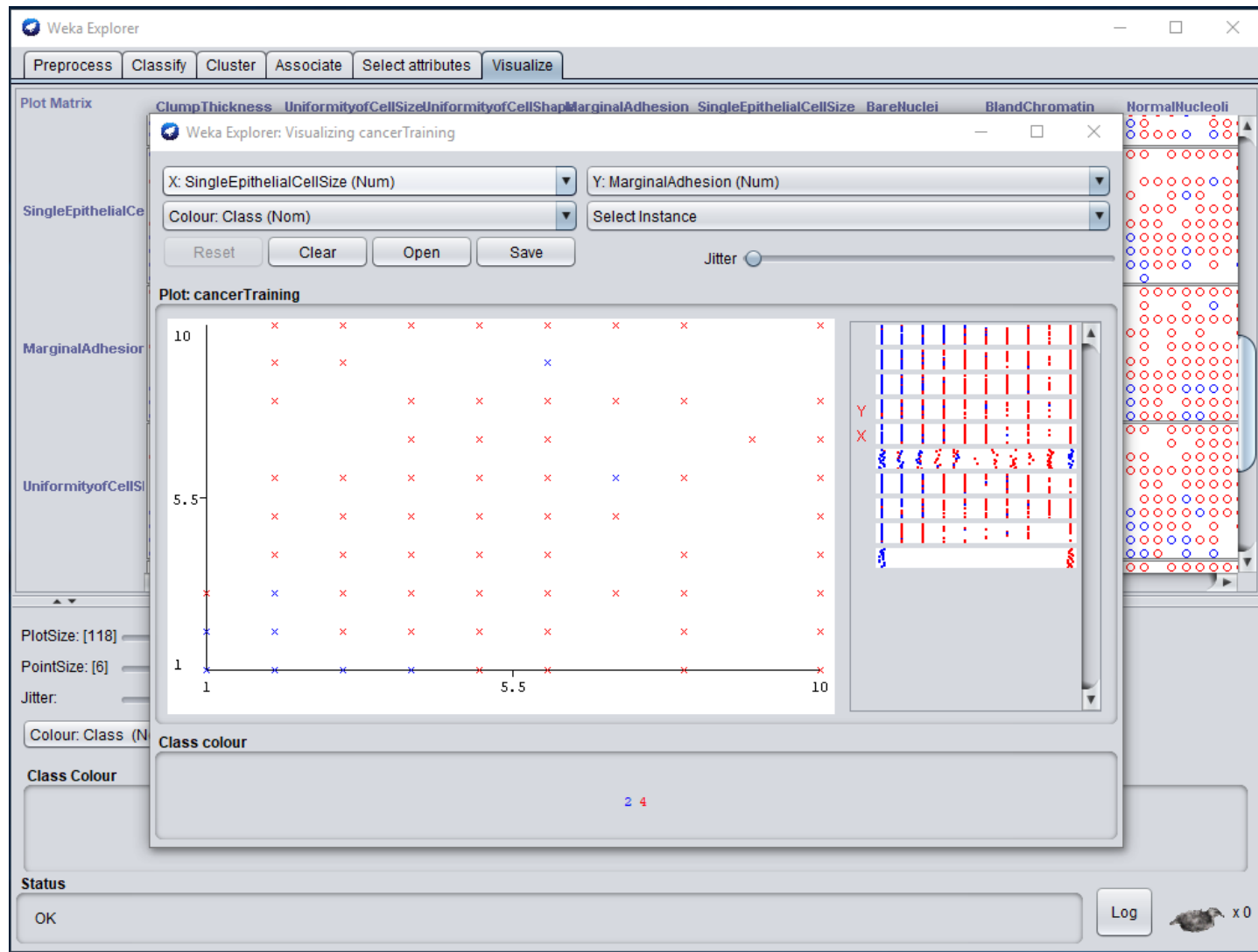
  - WEKA Data Mining Book:

    - Ian H. Witten and Eibe Frank, Data Mining: Practical Machine Learning Tools and Techniques (Second Edition)

  - WEKA Wiki:

    - http://weka.sourceforge.net/wiki/index.php/Main_Page

  - Others:

    - Jiawei Han and Micheline Kamber, Data Mining: Concepts and Techniques, 2nd ed