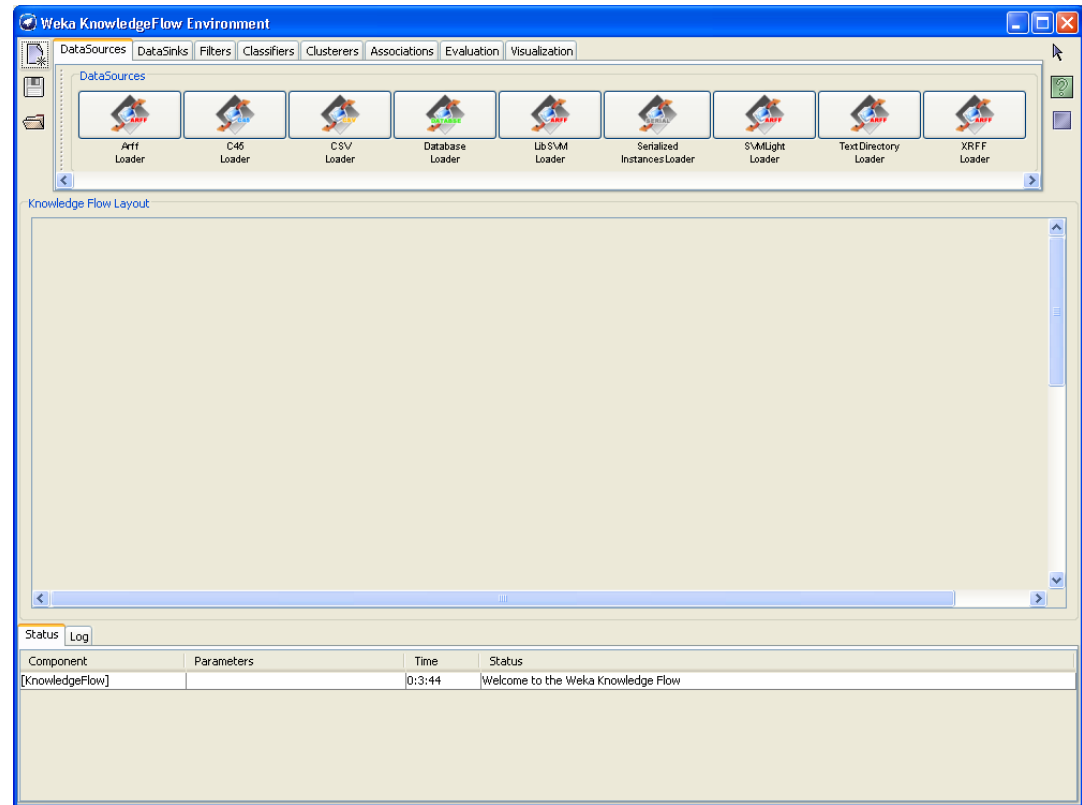# WEKA
# Waikato Environment for Knowledge Analysis

# Performing Classification Experiments

# Prof. Pietro Ducange

# The Knowledge Flow Interface

- It provides an alternative to the Explorer interface

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

# Knowledge Flow example (1)

Setting up a flow to load an ARFF file and perform a cross-validation using J48

- Create a source of data (DataSources tab - ARFFLoader)

- Connect it to a ARFF file (right click over the ARFFLoader icon - Configure)

- Specify which attribute is the class (Evaluation tab – ClassAssigner)

- Connect the ArffLoader to the ClassAssigner  (right click over the ArffLoader, select the dataSet under Connections and link with the ClassAssigner component with a left click

- Specify which column is the class (right click over the ClassAssigner - choose Configure)

- Add a CrossValidationFoldMaker component (Evaluation)

- Connect the ClassAssigner to the CrossValidationFoldMaker (right click over ClassAssigner, select dataSet, left click over CrossValidationFoldMaker

# Knowledge Flow example (2)

- Select the J48 component (classifiers tab)

- Connect the CrossValidationFoldMaker to J48 TWICE (right click over CrossValidationFoldMaker, first choose trainingSet and then testSet)

- Select ClassifierPerformanceEvaluator component (Evaluation tab)

- Connect J48 to this component (right click over J48, select batchClassifier left click over by ClassifierPerformanceEvaluator

- Select TextViewer component (Visualization tab)

- Connect the ClassifierPerformanceEvaluator to the TextViewer (select the text entry from the pop-up menu for ClassifierPerformanceEvaluator)

- Select GraphViewer component (Vizualization tab) and link to J48 (select the graph entry from the pop-up menu for J48)

- Start the flow (select start loading from the pop-up menu for the loader)

# Knowledge Flow example (3)

*Università di Pisa*

# Knowledge Flow example (4)



- Select show results from the pop-up menu for the graph viewer



- Select show results from the pop-up menu for the text viewer

6

# Knowledge Flow: attribute selection

# Knowledge Flow: attribute selection

Select show results from the pop-up menu for the text viewer connected to the Attribute Selection Block



**For each fold we can extracted the actual filtered test set!!!**

# Knowledge Flow: metaclassification

# Knowledge Flow: attribute selection

Select show results from the pop-up menu for the text viewer
connected to the Meta Classifier Block



**Text Viewer**

**Result list**

17:36:28 – Model: AttributeSelectedClassifier (fold 1)
17:36:28 – Model: AttributeSelectedClassifier (fold 2)
17:36:28 – Model: AttributeSelectedClassifier (fold 3)
17:36:28 – Model: AttributeSelectedClassifier (fold 4)
17:36:28 – Model: AttributeSelectedClassifier (fold 5)
17:36:28 – Model: AttributeSelectedClassifier (fold 6)
17:36:28 – Model: AttributeSelectedClassifier (fold 7)
17:36:28 – Model: AttributeSelectedClassifier (fold 8)
17:36:28 – Model: AttributeSelectedClassifier (fold 9)
17:36:28 – Model: AttributeSelectedClassifier (fold 10)

**Text**

```
Scheme:    AttributeSelectedClassifier
Relation: pima_diabetes
Training Fold: 10

AttributeSelectedClassifier:


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

Search Method:
        Best first.
        Start set: no attributes
        Search direction: forward
        Stale search after 5 node expansions
        Total number of subsets evaluated: 37
        Merit of best subset found:    0.162

Attribute Subset Evaluator (supervised, Class (nominal): 9 class):
        CFS Subset Evaluator
        Including locally predictive attributes

Selected attributes: 2,6,8 : 3
                     plas
                     mass
                     age


Header of reduced data:
@relation 'pima_diabetes-weka.filters.unsupervised.attribute.Remove-V-R2,6,8-9'

@attribute plas numeric
@attribute mass numeric
@attribute age numeric
@attribute class {tested_negative,tested_positive}

@data


Classifier Model
IB1 instance-based classifier
using 1 nearest neighbour(s) for classification
```

***For each fold we can extracted the actual model along with the selected features!***

# The Experimenter

- A robust experimental part involves running several learning schemes on different datasets.

- The Experimenter interface enables us to set-up large scale experiments.

- The user can create an experiment that runs several schemes against a series of datasets and then analyze the results to determine if one of the schemes is (statistically) better than the other schemes.

# Simple setup

**Experiment type:**

- Cross-validation (default), Train/Test Percentage Split (data randomized or order preserved)
- Number of folds
- Classification/Regression

**Iteration control**

- Set the number of repetition and change the order of iterations

**Datasets**

**Algorithms**

# The Analyze panel

The number of result lines available

**Weka Experiment Environment**

Setup | Run | Analyse

**Source**

Got 1200 results

Type of results to load:
→ from the current experiment
→ from an earlier experiment file
→ from the database

File... | Database... | Experiment

**Configure test**

| | |
|---|---|
| Testing with | Paired T–Test... ⇕ |
| Row | Select |
| Column | Select |
| Comparison field | Percent_correct |
| Significance | 0.05 |
| Sorting (asc.) by | <default> ⇕ |
| Test base | Select |
| Displayed Columns | Select |
| Show std. deviations | ☐ |
| Output Format | Select |

Perform test | Save output

**Result list**

17:06:17 – Available resultsets
17:06:39 – Percent_correct – trees.J48 '–C 0.25 –
17:11:16 – Percent_correct – trees.RandomForest

**Test output**

```
Tester:       weka.experiment.PairedCorrectedTTester
Analysing:    Percent_correct
Datasets:     4
Resultsets:   3
Confidence:   0.05 (two tailed)
Sorted by:    -
Date:         16/11/15 17.06


Dataset                      (1) trees.J4 | (2) trees (3) rules
---------------------------------------------------------------
pima_diabetes        (100)    74.49         74.44       75.18
Glass                (100)    67.63         76.16 v     66.78
ionosphere           (100)    89.74         93.11 v     89.16
iris                 (100)    94.73         94.27       93.93
---------------------------------------------------------------
                             (v/ /*) |     (2/2/0)    (0/4/0)

Key:
(1) trees.J48  –C 0.25 –M 2  –217733168393644444
(2) trees.RandomForest '–I 10 –K 0 –S 1' 4216839470751428698
(3) rules.JRip '–F 3 –N 2.0 –O 2 –S 1' –6589312996832147161
```

Type of comparison

Significance level

How to perform and show the results of the test

# The paired T-test results respect to a control algorithm (C4.5)



Weka Experiment Environment

Setup | Run | Analyse

**Source**

Got 1200results          File...   Database...   Experiment

**Configure test**

Testing with   Paired T–Test...

Row   Select

Column   Select

Comparison field   Percent_correct

Significance   0.05

Sorting (asc.) by   <default>

Test base   Select

Displayed Columns   Select

Show std. deviations   ☐

Output Format   Select

Perform test    Save output

**Result list**

17:06:17 – Available resultsets
17:06:39 – Percent_correct – trees.J48 '–C 0.25 –
17:11:16 – Percent_correct – trees.RandomForest

**Test output**

```
Tester:       weka.experiment.PairedCorrectedTTester
Analysing:    Percent_correct
Datasets:     4
Resultsets:   3
Confidence:   0.05 (two tailed)
Sorted by:    -
Date:         16/11/15 17.06


Dataset                    (1) trees.J4 | (2) trees (3) rules
-------------------------------------------------------------
pima_diabetes              (100)   74.49 |   74.44     75.18
Glass                      (100)   67.63 |   76.16 v   66.78
ionosphere                 (100)   89.74 |   93.11 v   89.16
iris                       (100)   94.73 |   94.27     93.93
-------------------------------------------------------------
                                   (v/ /*) | (2/2/0)   (0/4/0)


Key:
(1) trees.J48 '–C 0.25 –M 2' –217733168393644444
(2) trees.RandomForest '–I 10 –K 0 –S 1' 4216839470751428698
(3) rules.JRip '–F 3 –N 2.0 –O 2 –S 1' –6589312996832147161
```

v→ the results are statistically better than the control algorithm
*→ the results are statistically worse than the control algorithm
(x/y/z) → counts of the number of times
the scheme was better than (x), the same as (y), or worse than
(z)  the control algorithm

14

# The paired T-test results respect to a control algorithm (Random Forest)

```
Weka Experiment Environment

                    Setup    Run    Analyse

Source

Got 1200results                          File...    Database...    Experiment

Configure test                          Test output

Testing with   Paired T-Test...  ⬍     Tester:      weka.experiment.PairedCorrectedTTester
                                        Analysing:   Percent_correct
Row            Select                   Datasets:    4
                                        Resultsets: 3
Column         Select                   Confidence: 0.05 (two tailed)
                                        Sorted by:   -
                                        Date:        16/11/15 17.11
Comparison field  Percent_correct  ⬍

Significance   0.05                     Dataset                    (2) trees.Ra | (1) trees (3) rules
                                        ------------------------------------------------------------
Sorting (asc.) by  <default>     ⬍      pima_diabetes             (100)    74.44 |    74.49       75.18
                                        Glass                     (100)    76.16 |    67.63 *     66.78 *
Test base      Select                   ionosphere                (100)    93.11 |    89.74 *     89.16 *
                                        iris                      (100)    94.27 |    94.73       93.93
Displayed Columns  Select               ------------------------------------------------------------
                                                                   (v/ /*) |    (0/2/2)     (0/2/2)
Show std. deviations  ☐
                                        Key:
Output Format  Select                   (1) trees.J48 '-C 0.25 -M 2' -217733168393644444
                                        (2) trees.RandomForest '-I 10 -K 0 -S 1' 4216839470751428698
                                        (3) rules.JRip '-F 3 -N 2.0 -O 2 -S 1' -6589312996832147161
   Perform test        Save output

Result list

17:06:17 - Available resultsets
17:06:39 - Percent_correct - trees.J48 '-C 0.25 -I
17:11:16 - Percent_correct - trees.RandomForest
```

# Exercise

- Load the ionosphere dataset and prepare a 5 fold cross validation

- Perform the classification by using the three different classifiers and identify the most performing one

- Once selected the best classifier, perform the classification by using a metaclassifier with three different attribute selection methods

- Which is the best attribute selection method?

- Which are the most relevant selected attributes?