

More Data Mining with Weka

Class 1 – Lesson 1

Introduction

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

More Data Mining with Weka

**... a practical course on how to
use advanced facilities of Weka for data mining
(but not programming, just the interactive interfaces)**

... follows on from *Data Mining with Weka*

... will pick up some basic principles along the way

Ian H. Witten

University of Waikato, New Zealand

More Data Mining with Weka

❖ **This course assumes that you know about**

- *What data mining is and why it's useful*
- *The “simplicity-first” paradigm*
- *Installing Weka and using the Explorer interface*
- *Some popular classifier algorithms and filter methods*
- *Using classifiers and filters in Weka ...
and how to find out more about them*
- *Evaluating the result, including training/testing pitfalls*
- *Interpret Weka's output and visualizing your data set*
- *The overall data mining process*

❖ **See *Data Mining with Weka***

❖ **(Refresher: see videos on YouTube WekaMOOC channel)**

More Data Mining with Weka

❖ **As you know, a Weka is**

- *a bird found only in New Zealand?*
- **Data mining workbench:**
Waikato Environment for Knowledge Analysis

Machine learning algorithms for data mining tasks

- 100+ algorithms for classification
- 75 for data preprocessing
- 25 to assist with feature selection
- 20 for clustering, finding association rules, etc

More Data Mining with Weka

What will you learn?

- ❖ *Experimenter, Knowledge Flow interface, Command Line interfaces*
- ❖ *Dealing with “big data”*
- ❖ *Text mining*
- ❖ *Supervised and unsupervised filters*
- ❖ *All about discretization, and sampling*
- ❖ *Attribute selection methods*
- ❖ *Meta-classifiers for attribute selection and filtering*
- ❖ *All about classification rules: rules vs. trees, producing rules*
- ❖ *Association rules and clustering*
- ❖ *Cost-sensitive evaluation and classification*

Use Weka on your own data ... and understand what you're doing!

Class 1: Exploring Weka's interfaces, and working with big data

- ❖ Experimenter interface
- ❖ Using the Experimenter to compare classifiers
- ❖ Knowledge Flow interface
- ❖ Simple Command Line interface
- ❖ Working with big data
 - *Explorer: 1 million instances, 25 attributes*
 - *Command line interface: effectively unlimited*
 - *in the Activity you will process a multi-million-instance dataset*

Course organization

**Class 1 Exploring Weka's interfaces;
working with big data**

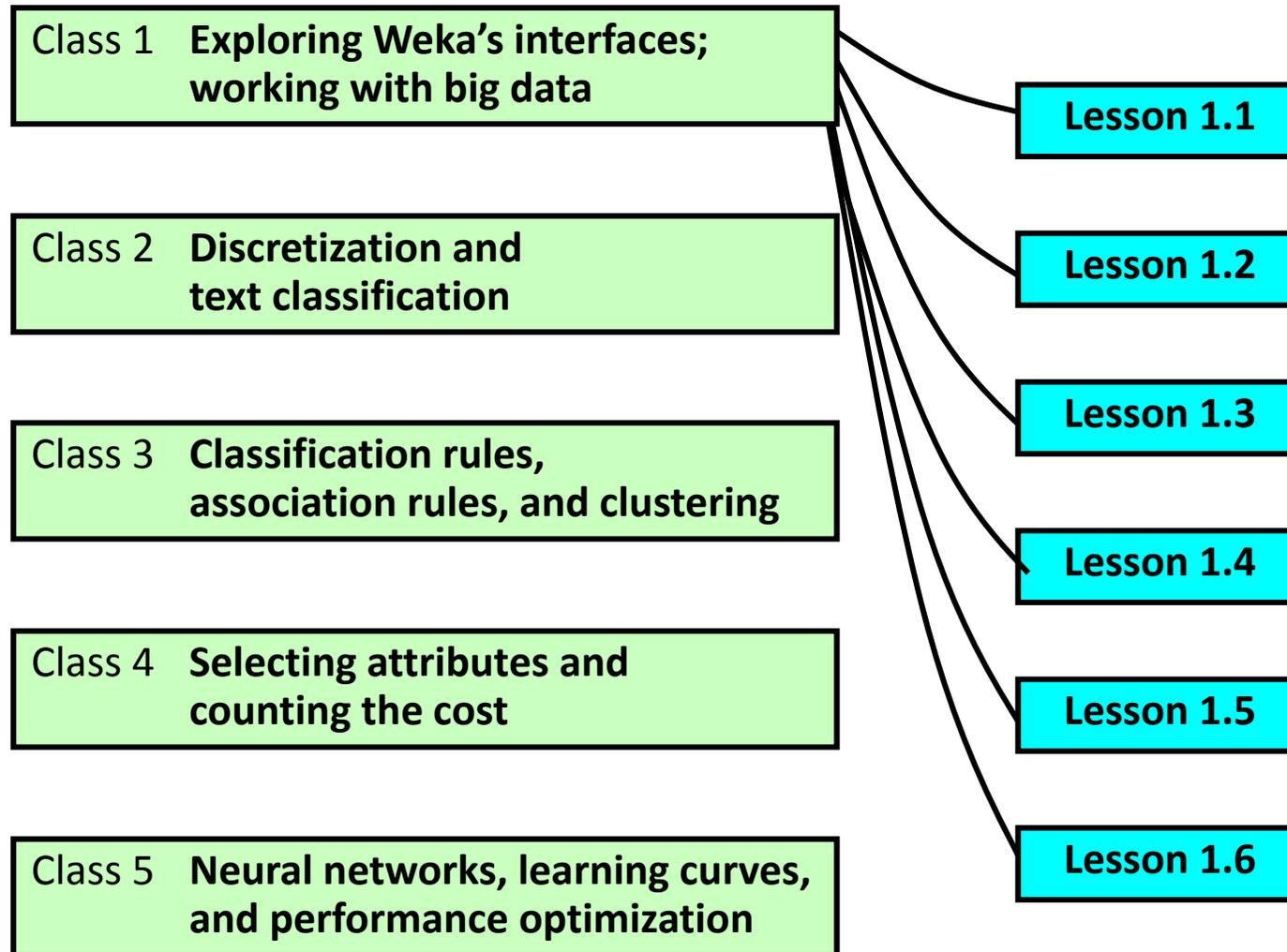
**Class 2 Discretization and
text classification**

**Class 3 Classification rules,
association rules, and clustering**

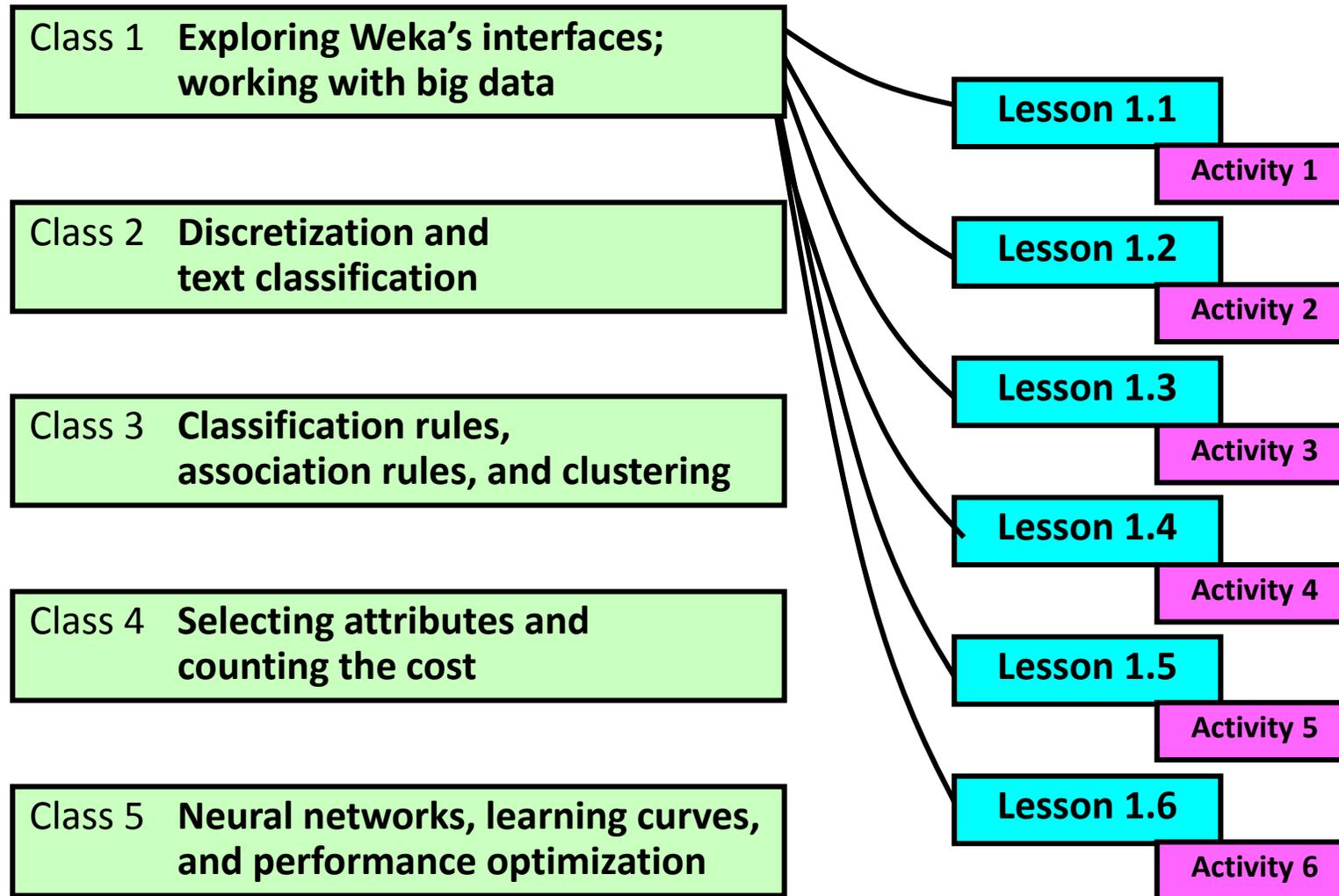
**Class 4 Selecting attributes and
counting the cost**

**Class 5 Neural networks, learning curves,
and performance optimization**

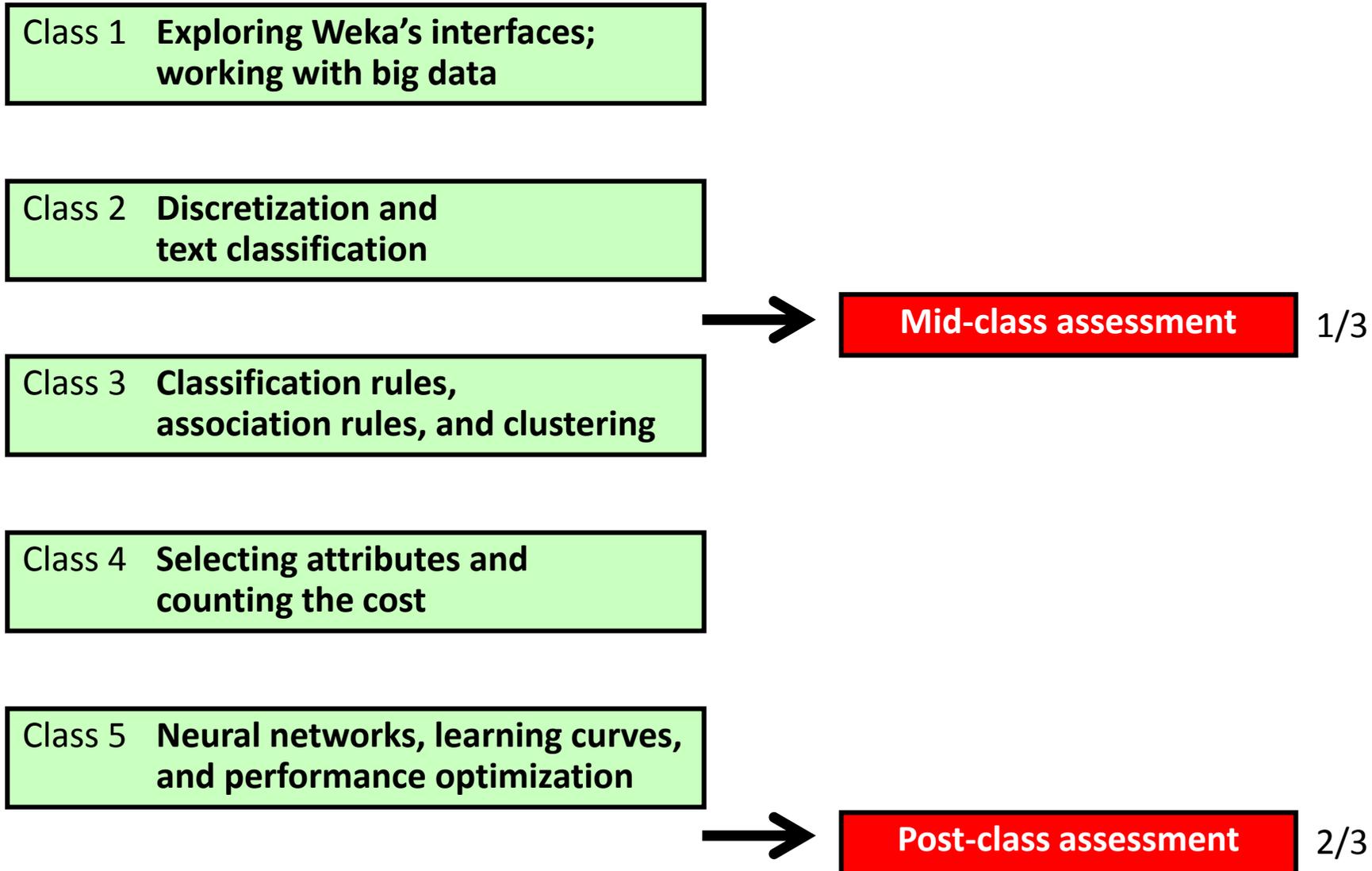
Course organization



Course organization



Course organization



Download Weka now!

Download from

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

for Windows, Mac, Linux

Weka 3.6.11

the latest stable version of Weka

includes datasets for the course

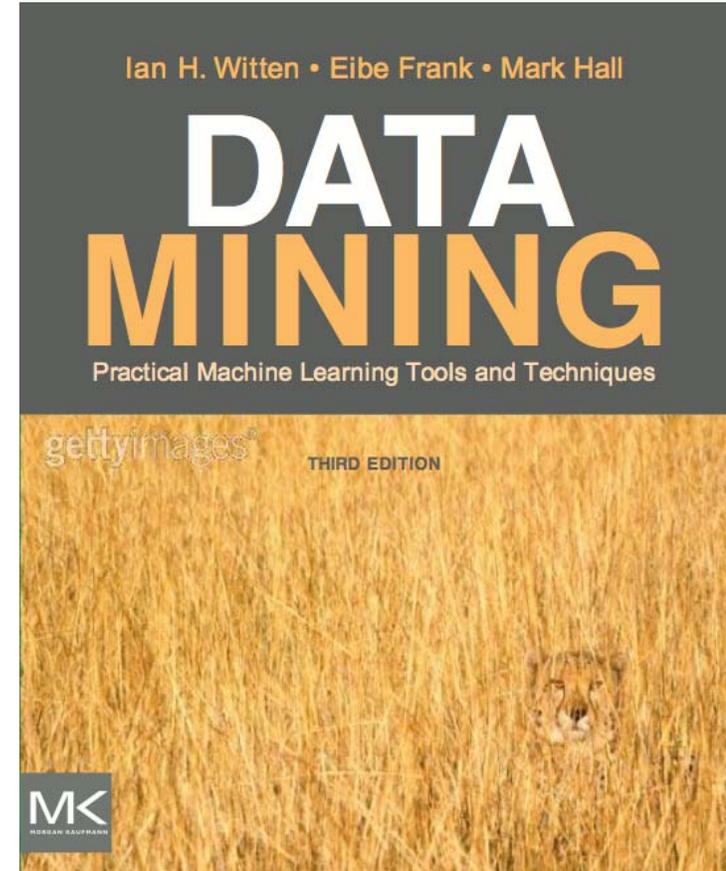
it's important to get the right version!

Textbook

This textbook discusses data mining, and Weka, in depth:

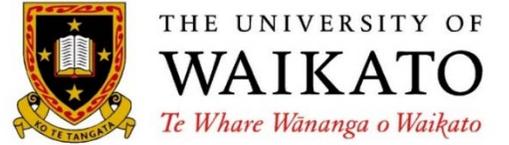
Data Mining: Practical machine learning tools and techniques,
by Ian H. Witten, Eibe Frank and
Mark A. Hall. Morgan Kaufmann, 2011

The publisher has made available parts relevant to this course in ebook format.









More Data Mining with Weka

Class 1 – Lesson 2

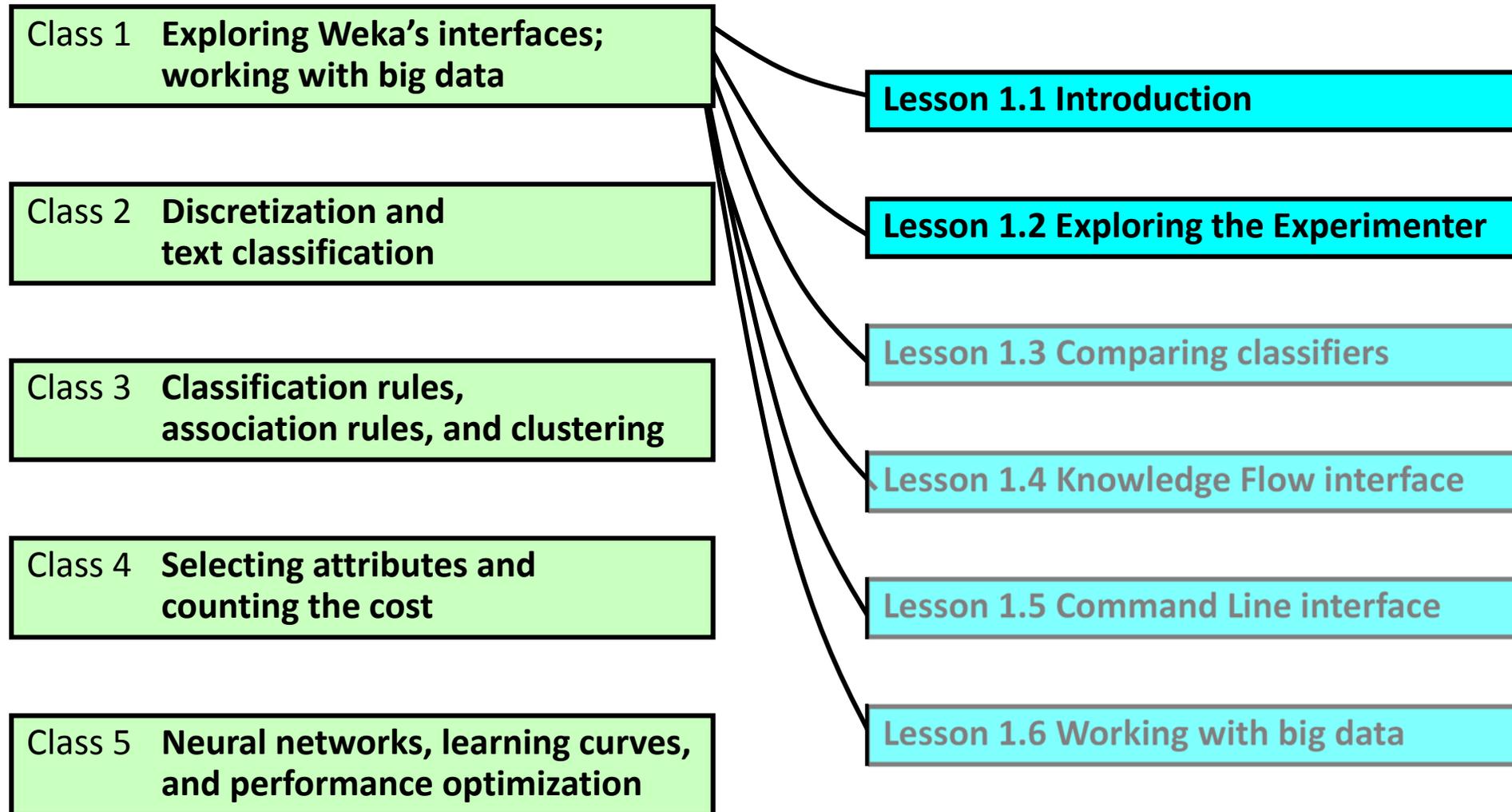
Exploring the Experimenter

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

Lesson 1.2: Exploring the Experimenter



Lesson 1.2: Exploring the Experimenter



Trying out classifiers/filters

Performance comparisons

Graphical interface

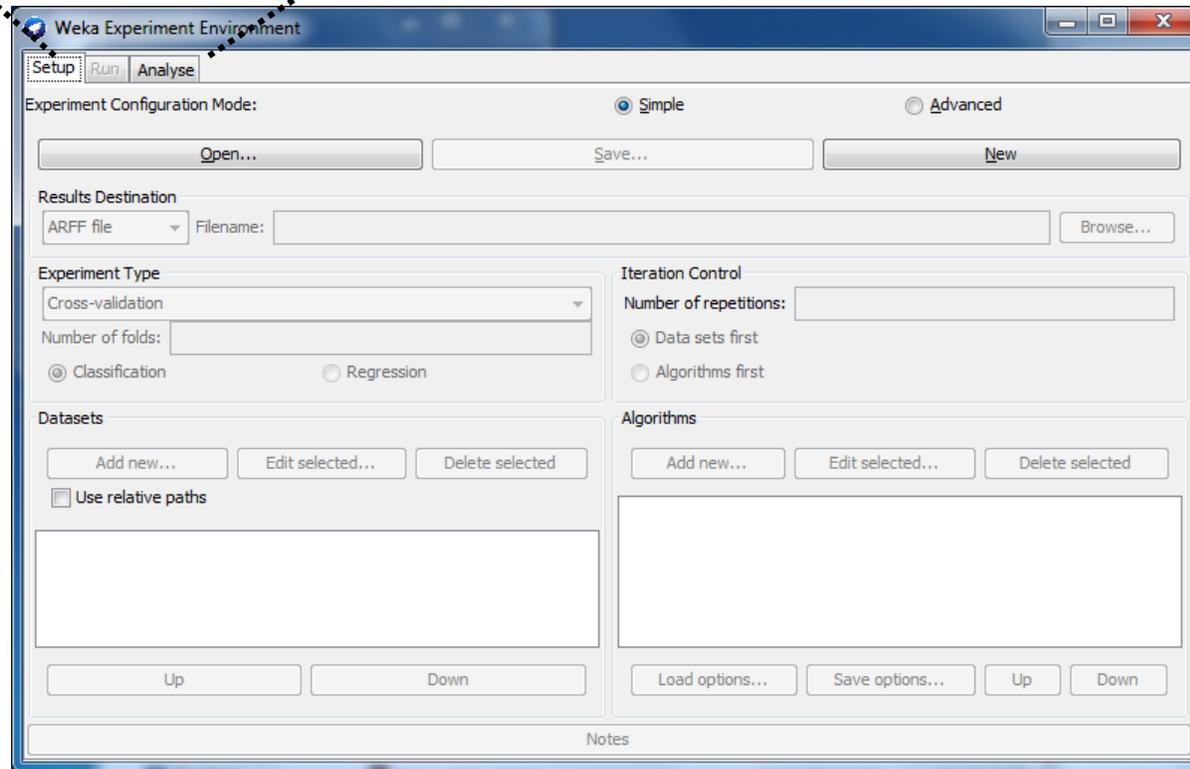
Command-line interface

Lesson 1.2: Exploring the Experimenter

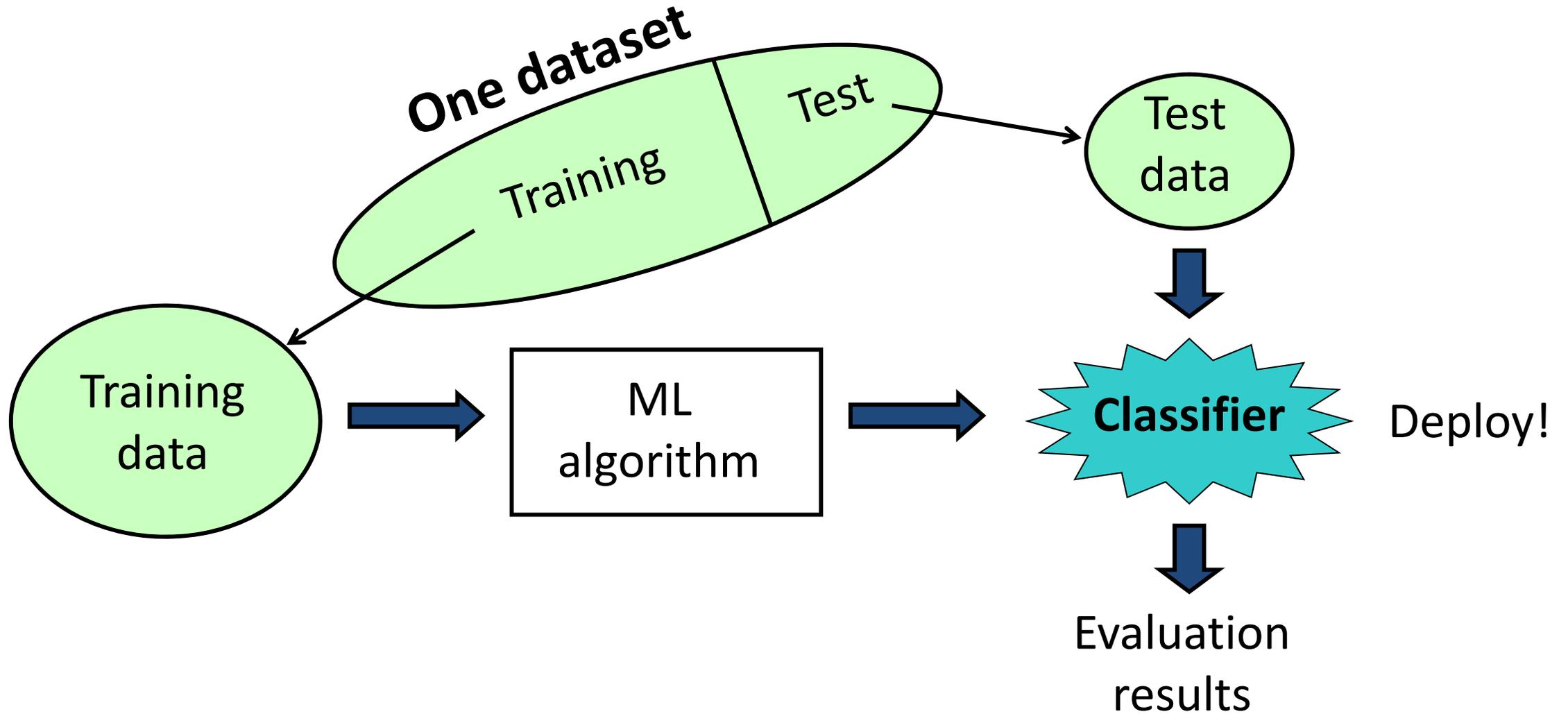
Use the Experimenter for ...

- ❖ determining mean and standard deviation performance of a classification algorithm on a dataset
... or several algorithms on several datasets
- ❖ Is one classifier better than another on a particular dataset?
... and is the difference statistically significant?
- ❖ Is one parameter setting for an algorithm better than another?
- ❖ The result of such tests can be expressed as an ARFF file
- ❖ Computation may take days or weeks
... and can be distributed over several computers

Lesson 1.2: Exploring the Experimenter



Lesson 1.2: Exploring the Experimenter

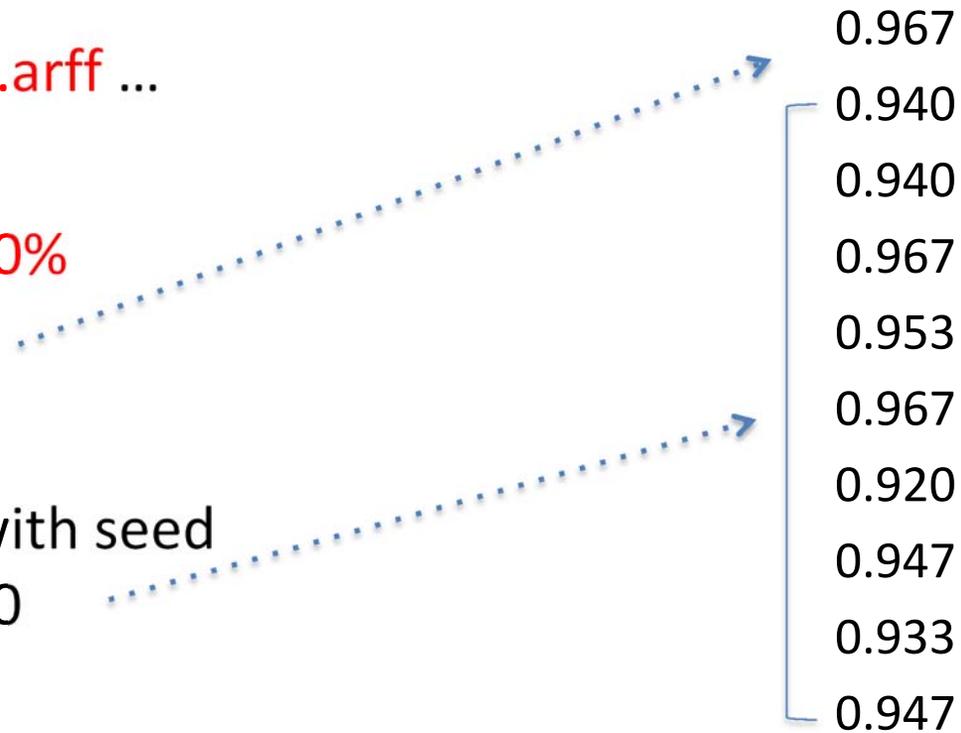


Basic assumption: training and test sets produced by independent sampling from an infinite population

Lesson 1.2: Exploring the Experimenter

Evaluate J48 on segment-challenge (*Data Mining with Weka*, Lesson 2.3)

- ❖ With **segment-challenge.arff** ...
- ❖ and J48 (**trees>J48**)
- ❖ Set **percentage split** to **90%**
- ❖ Run it: 96.7% accuracy
- ❖ Repeat
- ❖ [**More options**] Repeat with seed
2, 3, 4, 5, 6, 7, 8, 9 10



Lesson 1.2: Exploring the Experimenter

Evaluate J48 on segment-challenge (*Data Mining with Weka, Lesson 2.3*)

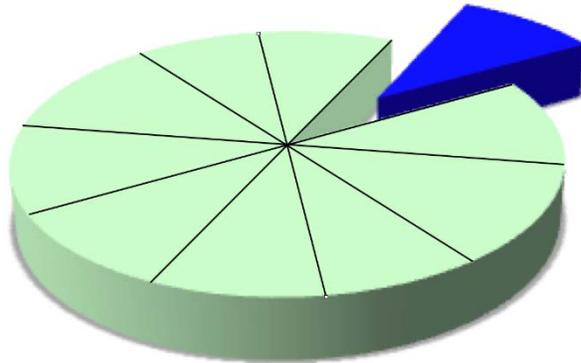
Sample mean	$\bar{x} = \frac{\sum x_i}{n}$	0.967
		0.940
		0.940
		0.967
Variance	$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1}$	0.953
		0.967
		0.920
Standard deviation	σ	0.947
		0.933
		0.947

$$\bar{x} = 0.949, \sigma = 0.018$$

Lesson 1.2: Exploring the Experimenter

10-fold cross-validation (*Data Mining with Weka, Lesson 2.5*)

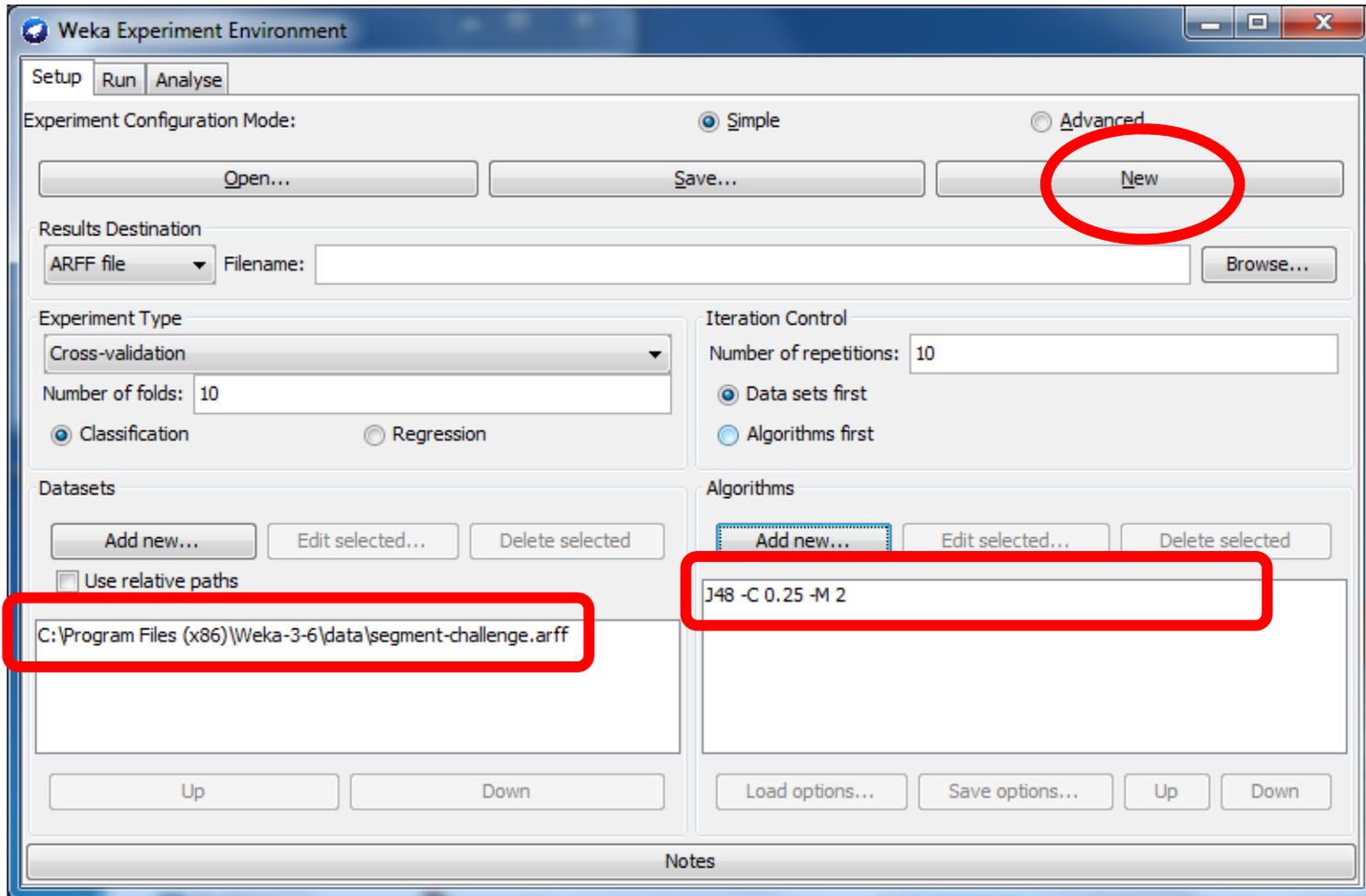
- ❖ Divide dataset into 10 parts (folds)
- ❖ Hold out each part in turn
- ❖ Average the results
- ❖ Each data point used once for testing, 9 times for training



Stratified cross-validation

- ❖ Ensure that each fold has the right proportion of each class value

Lesson 1.2: Exploring the Experimenter



Setup panel

- ❖ click *New*
- ❖ note defaults
 - 10-fold cross-validation, repeat 10 times
- ❖ under *Datasets*, click *Add new*, open [segment-challenge.arff](#)
- ❖ under *Algorithms*, click *Add new*, open [trees>J48](#)

Run panel

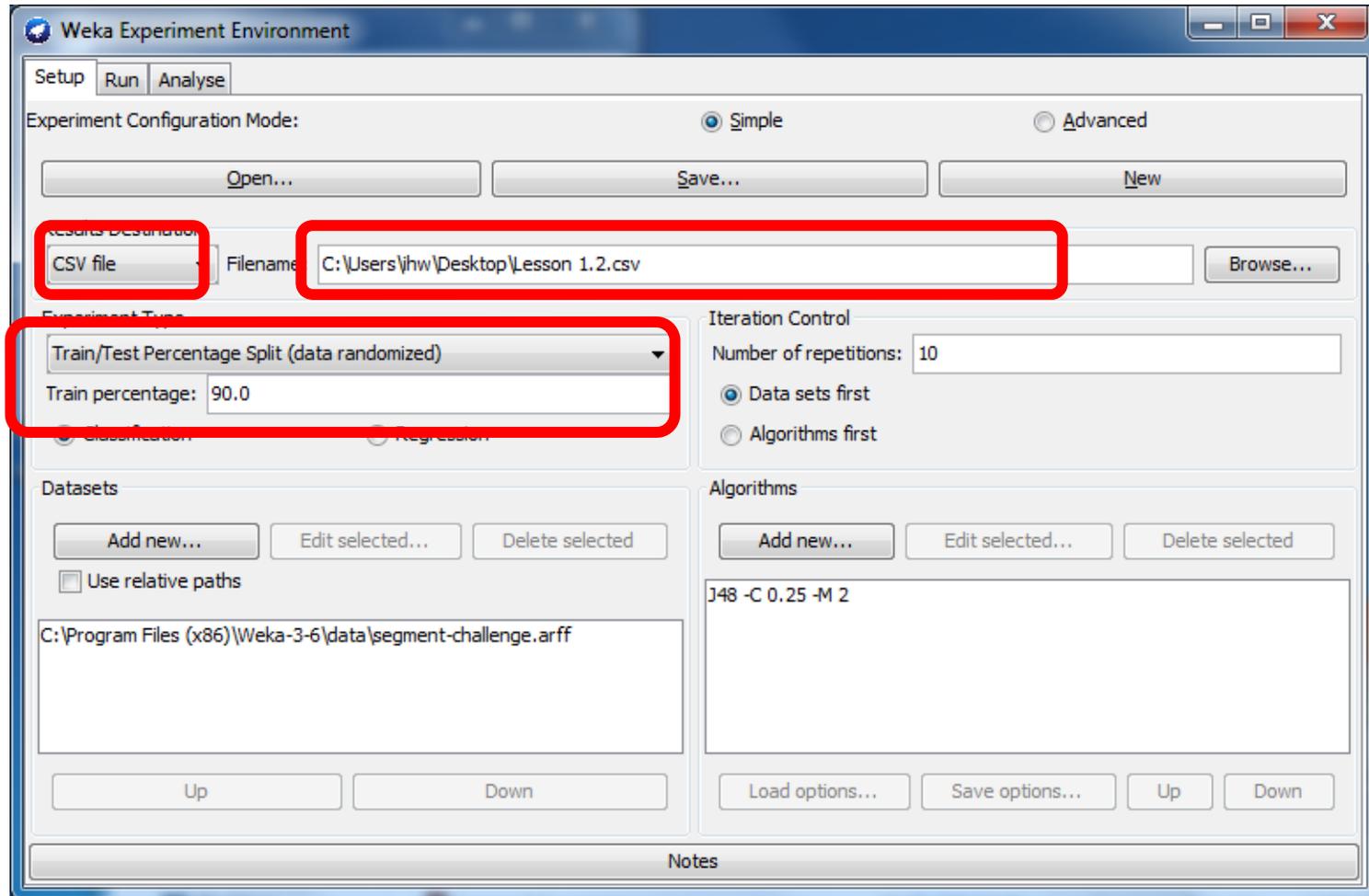
- ❖ click *Start*

Analyse panel

- ❖ click *Experiment*
- ❖ Select *Show std. deviations*
- ❖ Click *Perform test*

$$\bar{x} = 95.71\%, \sigma = 1.85\%$$

Lesson 1.2: Exploring the Experimenter



To get detailed results

return to *Setup* panel

- ❖ select .csv file
- ❖ enter filename for results
- ❖ *Train/Test Split; 90%*

Lesson 1.2: Exploring the Experimenter

Re-run cross-validation experiment

- ❖ Open results spreadsheet

Lesson 1.2: Exploring the Experimenter

Setup panel

- ❖ Save/Load an experiment
- ❖ Save the results in Arff file ... or in a database
- ❖ Preserve order in Train/Test split (can't do repetitions)
- ❖ Use several datasets, and several classifiers
- ❖ Advanced mode

Run panel

Analyse panel

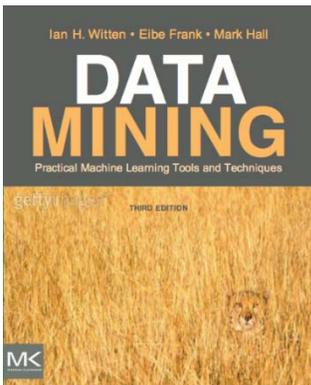
- ❖ Load results from .csv or Arff file ... or from a database
- ❖ Many options

Lesson 1.2: Exploring the Experimenter

- ❖ Open Experimenter
- ❖ Setup, Run, Analyse panels
- ❖ Evaluate one classifier on one dataset
 - ... using cross-validation, repeated 10 times
 - ... using percentage split, repeated 10 times
- ❖ Examine spreadsheet output
- ❖ Analyse panel to get mean and standard deviation
- ❖ Other options on Setup and Run panels

Course text

- ❖ Chapter 13 *The Experimenter*





More Data Mining with Weka

Class 1 – Lesson 3

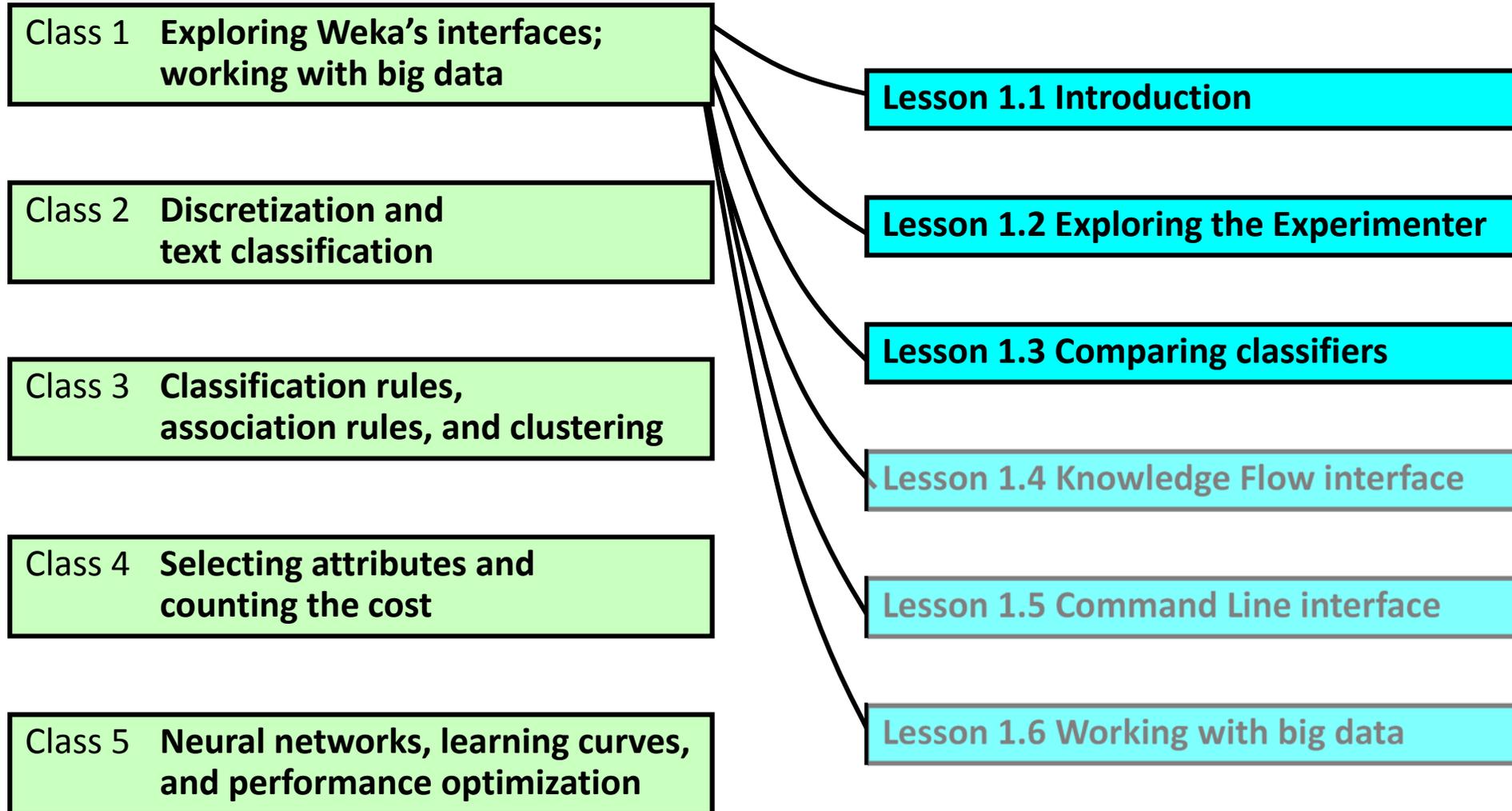
Comparing classifiers

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

Lesson 1.3: Comparing classifiers



Lesson 1.3: Comparing classifiers

Is J48 better than (a) ZeroR and (b) OneR on the Iris data?

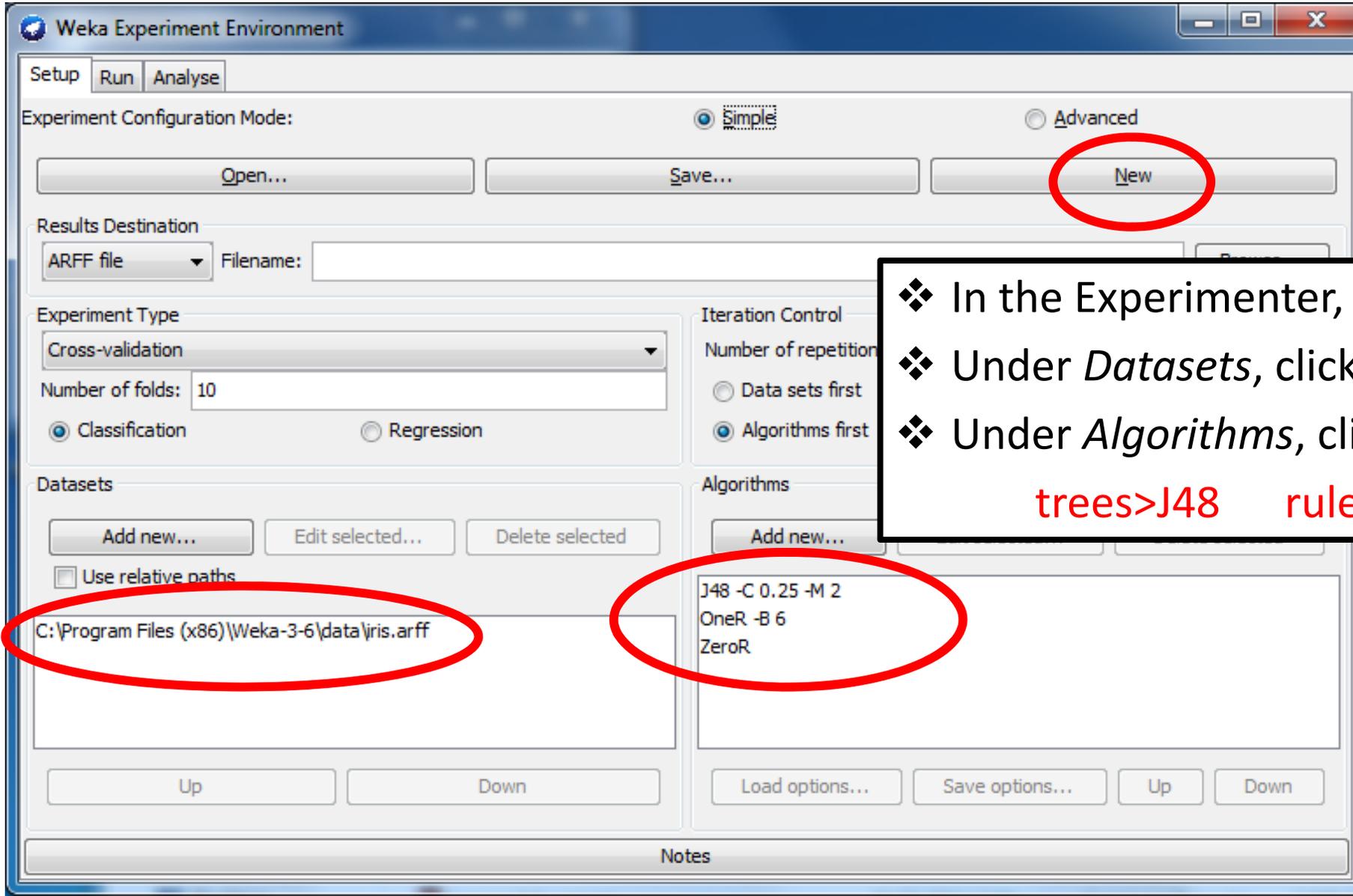
- ❖ In the Explorer, open `iris.arff`
- ❖ Using cross-validation, evaluate classification accuracy with ...

ZeroR (<code>rules>ZeroR</code>)	33%
OneR (<code>rules>OneR</code>)	92%
J48 (<code>trees>J48</code>)	96%

But how reliable is this?

What would happen if you used a different random number seed??

Lesson 1.3: Comparing classifiers



- ❖ In the Experimenter, click *New*
- ❖ Under *Datasets*, click *Add new*, open *iris.arff*
- ❖ Under *Algorithms*, click *Add new*, open
trees>J48 *rules>OneR* *rules>ZeroR*

Lesson 1.3: Comparing classifiers

The screenshot shows the Weka Experiment Environment window. The 'Analyse' tab is active. In the 'Configure test' section, the 'Testing with' dropdown is set to 'Paired T-Tester (corrected)'. The 'Significance' field is set to '0.05'. The 'Perform test' button is highlighted with a red circle. In the 'Test output' section, the 'Experiment' button is circled in red. Below the 'Test output' section, a table of results is displayed, which is also circled in red.

Source
Got 300 results

File... Database... Experiment

Configure test

Testing with: Paired T-Tester (corrected)

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

Test output

Tester: weka.experiment.PairedCorrected
Analysing: Percent_correct
Datasets: 1
Resultsets: 3
Confidence: 0.05 (two tailed)
Sorted by: -
Date: 4/03/14 10:41 AM

Dataset	(1) trees.J4	(2) rules	(3) rules
iris	(100) 94.73	92.53	33.33 *
	(v/ /*)	(0/1/0)	(0/0/1)

Key:
(1) trees.J48 '-C 0.25 -M 2' -217733168393644444
(2) rules.OneR '-B 6' -3459427003147861443
(3) rules.ZeroR '' 48055541465867954

- ❖ Switch to *Run*; click *Start*
- ❖ Switch to *Analyse*, click *Experiment*
click *Perform test*

Lesson 1.3: Comparing classifiers

Dataset	(1) trees.J48	(2) rules.OneR	(3) rules.ZeroR
iris	(100) 94.73	92.53	33.33 *
	(v/ /*)	(0/1/0)	(0/0/1)

Key:
(1) trees.J48
(2) rules.OneR
(3) rules.ZeroR

v significantly better
* significantly worse

- ❖ ZeroR (33.3%) is significantly worse than J48 (94.7%)
- ❖ Cannot be sure that OneR (92.5%) is significantly worse than J48
- ❖ ... at the 5% level of statistical significance
- ❖ J48 seems better than ZeroR: pretty sure (5% level) that this is not due to chance
- ❖ ... and better than OneR; but this may be due to chance (can't rule it out at 5% level)

Lesson 1.3: Comparing classifiers

Dataset	(1) trees.J4	(2) rules	(3) rules
iris	(100) 94.73	92.53	33.33 *
breast-cancer	(100) 74.28	66.91 *	70.30
german_credit	(100) 71.25	65.91 *	70.00
pima_diabetes	(100) 74.49	71.52	65.11 *
Glass	(100) 67.63	57.40 *	35.51 *
ionosphere	(100) 89.74	82.28 *	64.10 *
segment	(100) 95.71	64.35 *	15.73 *
	(v/ /*)	(0/2/5)	(0/2/5)

Key:

(1) trees.J48
(2) rules.OneR
(3) rules.ZeroR

J48 is significantly (5% level) better than

- ❖ both OneR and ZeroR on Glass, ionosphere, segment
- ❖ OneR on breast-cancer, german_credit
- ❖ ZeroR on iris, pima_diabetes

Lesson 1.3: Comparing classifiers

Dataset	(2) rules.On	(1) trees	(3) rules
iris	(100) 92.53	94.73	33.33 *
breast-cancer	(100) 66.91	74.28 v	70.30
german_credit	(100) 65.91	71.25 v	70.00 v
pima_diabetes	(100) 71.52	74.49	65.11 *
Glass	(100) 57.40	67.63 v	35.51 *
ionosphere	(100) 82.28	89.74 v	64.10 *
segment	(100) 64.35	95.71 v	15.73 *

(v/ /*) | (5/2/0) (1/1/5)

Key:

(1) trees.J48
(2) rules.OneR
(3) rules.ZeroR

Comparing OneR with ZeroR

Change “Test base” on Analyse panel

- ❖ significantly worse on german-credit
- ❖ about the same on breast-cancer
- ❖ significantly better on all the rest

Lesson 1.3: Comparing classifiers

The screenshot shows the Weka Experiment Environment interface. The 'Configure test' panel on the left is set to 'Paired T-Tester (correctness)'. The 'Row' and 'Column' buttons are highlighted with red boxes. The 'Comparison field' is 'Percent_correct', 'Significance' is 0.05, and 'Sorting (asc.) by' is '<default>'. The 'Perform test' button at the bottom left is also highlighted with a red box. The 'Test output' panel on the right shows the results of the test, including a table comparing three datasets: (1) iris, (2) breast-cancer, and (3) german_credit. The table shows accuracy percentages for three classifiers: trees.J48, rules.OneR, and rules.ZeroR. The 'Perform test' button is highlighted with a red box.

Source
Got 2100 results

File... Database... Experiment

Configure test

Testing with: Paired T-Tester (correctness)

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

Test output

Confidence: 0.05 (two tailed)

Sorted by: -

Date: 4/03/14 1:56 PM

Dataset	(1) iris	(2) breast-cancer	(3) german_credit
trees.J48	(100) 94.73	74.28 *	71.25 *
rules.OneR	(100) 92.53	66.91 *	65.91 *
rules.ZeroR	(100) 33.33	70.30 v	70.00 v

Key:

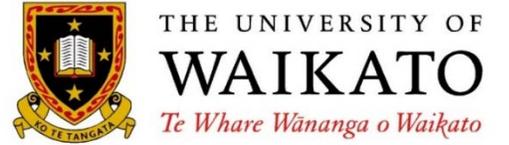
- (1) iris
- (2) breast-cancer
- (3) german_credit
- (4) pima_diabetes
- (5) Glass

❖ Row: select *Scheme* (not *Dataset*)

❖ Column: select *Dataset* (not *Scheme*)

Lesson 1.3: Comparing classifiers

- ❖ Statistical significance: the “null hypothesis”
Classifier A’s performance is the same as B’s
- ❖ The observed result is highly unlikely if the null hypothesis is true
“The null hypothesis can be rejected at the 5% level”
[of statistical significance]
“A performs significantly better than B at the 5% level”
- ❖ Can change the significance level (5% and 1% are common)
- ❖ Can change the comparison field (we have used % correct)
- ❖ Common to compare over a set of datasets
“On these datasets, method A has xx wins and yy losses over method B”
- ❖ Multiple comparison problem
if you make many tests, some will appear to be “significant” just by chance!



More Data Mining with Weka

Class 1 – Lesson 4

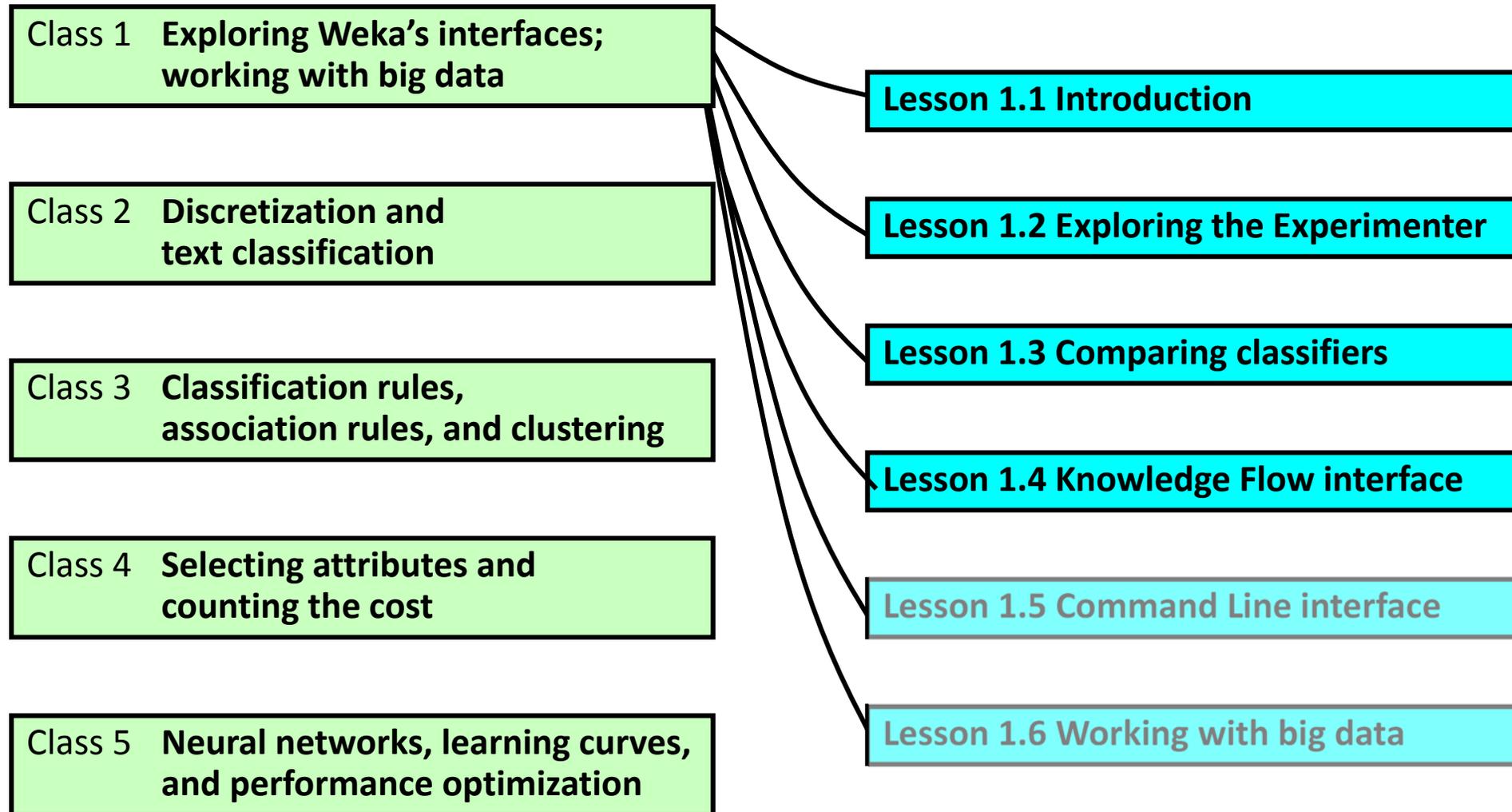
The Knowledge Flow interface

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

Lesson 1.4: The Knowledge Flow interface



Lesson 1.4: The Knowledge Flow interface

The Knowledge Flow interface is an alternative to the Explorer

- ❖ Lay out filters, classifiers, evaluators interactively on a 2D canvas
- ❖ Components include data sources, data sinks, evaluation, visualization
- ❖ Different kinds of connections between the components
 - *Instance or dataset*
 - *test set, training set*
 - *classifier*
 - *output, text or chart*
- ❖ Can work incrementally, on potentially infinite data streams
- ❖ Can look inside cross-validation at the individual models produced

Lesson 1.4: The Knowledge Flow interface

Load an ARFF file, choose J48, evaluate using cross-validation

	Toolbar
❖ Choose an <i>ArffLoader</i> ; Configure to set the file <i>iris.arff</i>	<i>DataSources</i>
❖ Connect up a <i>ClassAssigner</i> to select the class	<i>Evaluation</i>
❖ Connect the result to a <i>CrossValidationFoldMaker</i>	<i>Evaluation</i>
❖ Connect this to <i>J48</i>	<i>Classifiers</i>
❖ Make two connections, one for <i>trainingSet</i> and the other for <i>testSet</i>	
❖ Connect <i>J48</i> to <i>ClassifierPerformanceEvaluator</i>	<i>Evaluation</i>
❖ Connect this to a <i>TextViewer</i>	<i>Visualization</i>

Then run it! (*ArffLoader*: Start loading)

Lesson 1.4: The Knowledge Flow interface

The screenshot displays the Weka KnowledgeFlow Environment interface. At the top, there is a menu bar with tabs for Filters, Classifiers, Clusters, Associations, Evaluation, and Visualization. Below the menu bar is a toolbar with icons for various visualization tools: Data Visualizer, Scatter PlotMatrix, Attribute Summarizer, Model PerformanceChart, CostBenefit Analysis, Text Viewer, Graph Viewer, and Strip Chart.

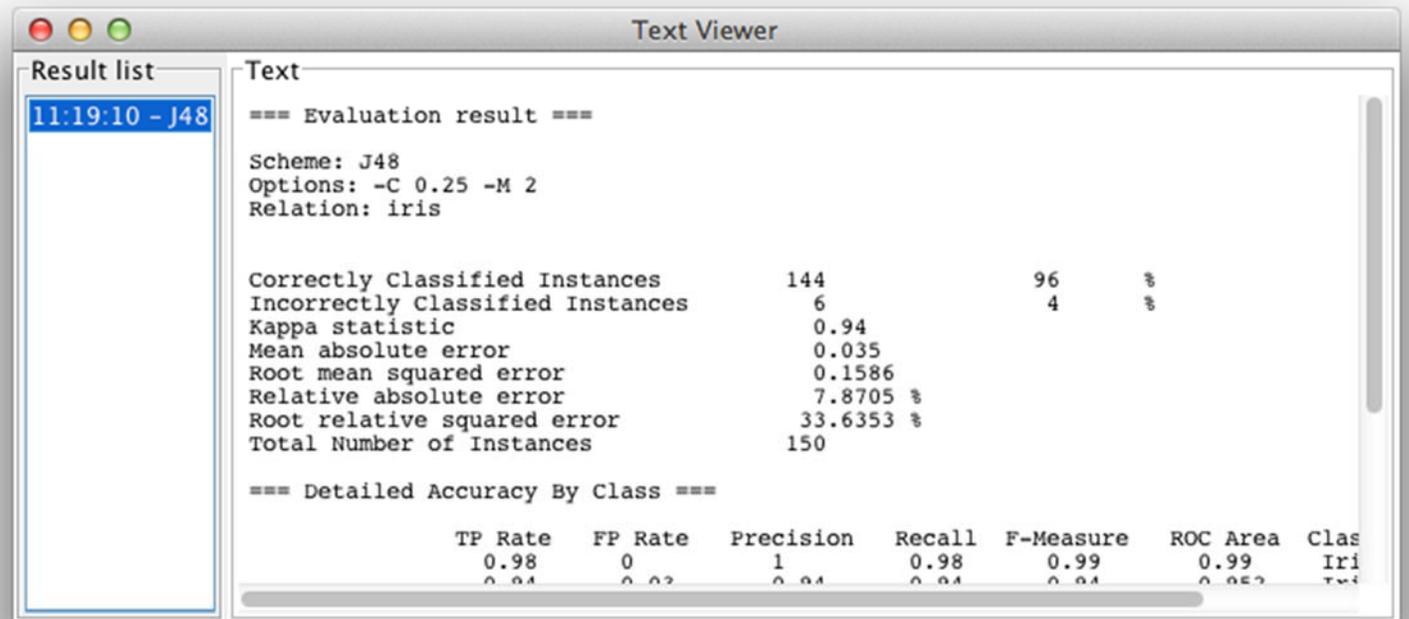
The main area is titled "Knowledge Flow Layout" and contains a workflow diagram. The workflow starts with an "ArffLoader" component, which outputs a "dataSet" to a "ClassAssigner" component. The "ClassAssigner" outputs a "dataSet" to a "CrossValidation FoldMaker" component. The "CrossValidation FoldMaker" outputs a "trainingSet" and a "testSet" to a "J48" component. The "J48" component outputs a "batchClassifier" to a "ClassifierPer..." component. The "ClassifierPer..." component outputs a "text" file to a "Text Viewer" component.

At the bottom of the interface, there is a "Status" and "Log" window. The log window shows the following entries:

Component	Parameters	Ti...	Status
[KnowledgeF...		0...	Welcome to the Weka Knowledge Flow
ArffLoader		-	Finished.
CrossValidat...		-	Finished.
J48	-C 0.25 -M 2	-	Finished.
ClassifierPer...		-	Finished.

Lesson 1.4: The Knowledge Flow interface

❖ *TextViewer*: Show results



```
Text Viewer

Result list
11:19:10 - J48

Text

=== Evaluation result ===

Scheme: J48
Options: -C 0.25 -M 2
Relation: iris

Correctly Classified Instances      144          96  %
Incorrectly Classified Instances    6            4  %
Kappa statistic                    0.94
Mean absolute error                 0.035
Root mean squared error             0.1586
Relative absolute error             7.8705 %
Root relative squared error        33.6353 %
Total Number of Instances          150

=== Detailed Accuracy By Class ===

      TP Rate  FP Rate  Precision  Recall  F-Measure  ROC Area  Clas
      0.98     0       1          0.98   0.99       0.99     Iri
      0.04     0.02    0.04       0.04   0.04       0.052     Tex
```

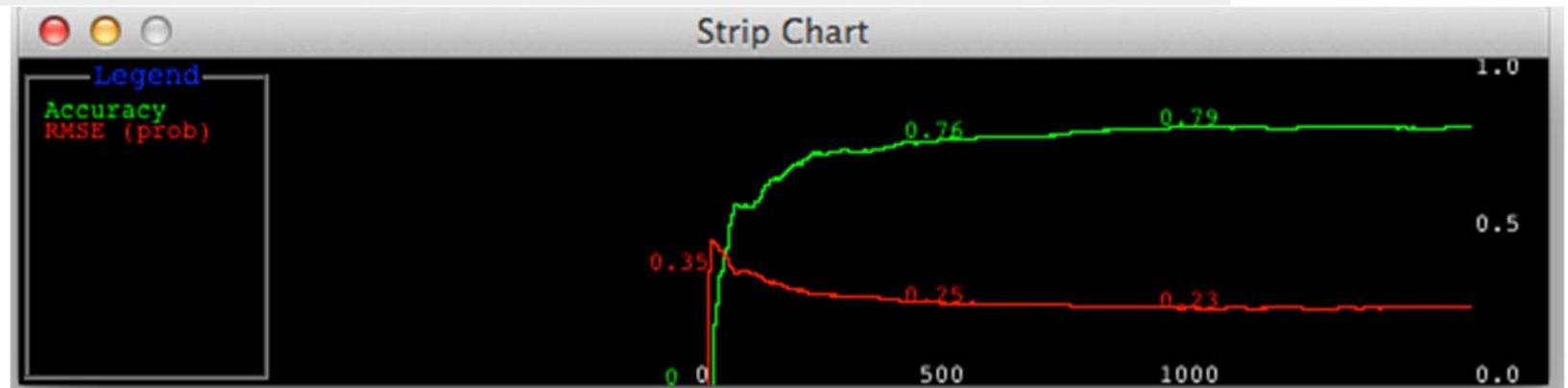
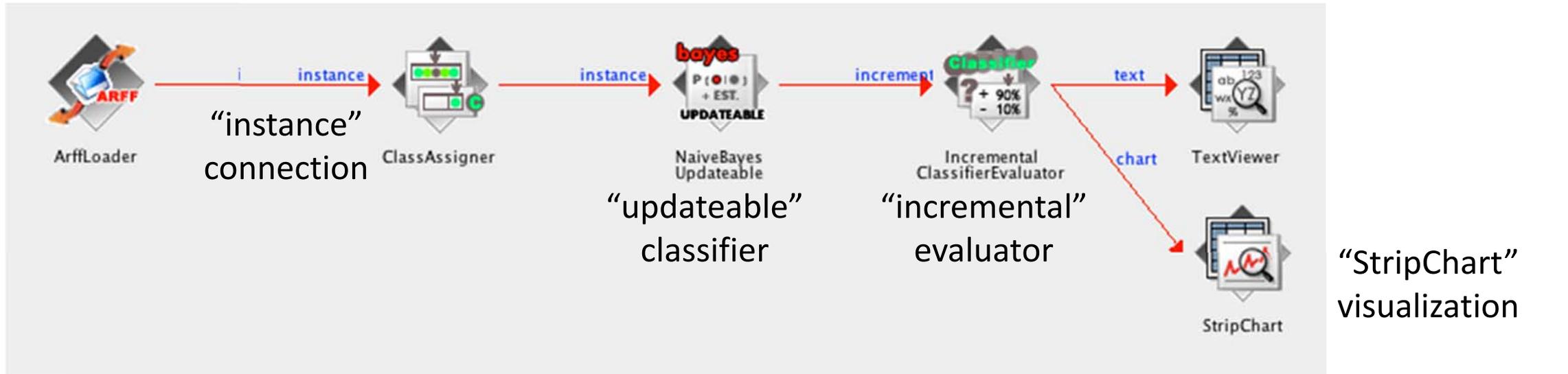
❖ Add a *ModelPerformanceChart*

❖ Connect the *visualizableError* output of *ClassifierPerformanceEvaluator* to it

❖ Show chart (need to run again)

Lesson 1.4: The Knowledge Flow interface

Working with stream data

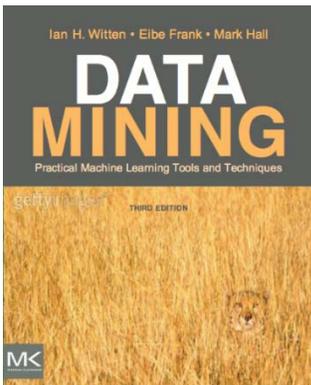


Lesson 1.4: The Knowledge Flow interface

- ❖ Panels broadly similar to the Explorer's, except
 - *DataSources* are separate from *Filters*
 - Write data or models to files using *DataSinks*
 - *Evaluation* is a separate panel
- ❖ Facilities broadly similar too, except
 - Can deal incrementally with potentially infinite datasets
 - Can look inside cross-validation at the models for individual folds
- ❖ Some people like graphical interfaces

Course text

- ❖ Chapter 12 *The Knowledge Flow Interface*





More Data Mining with Weka

Class 1 – Lesson 5

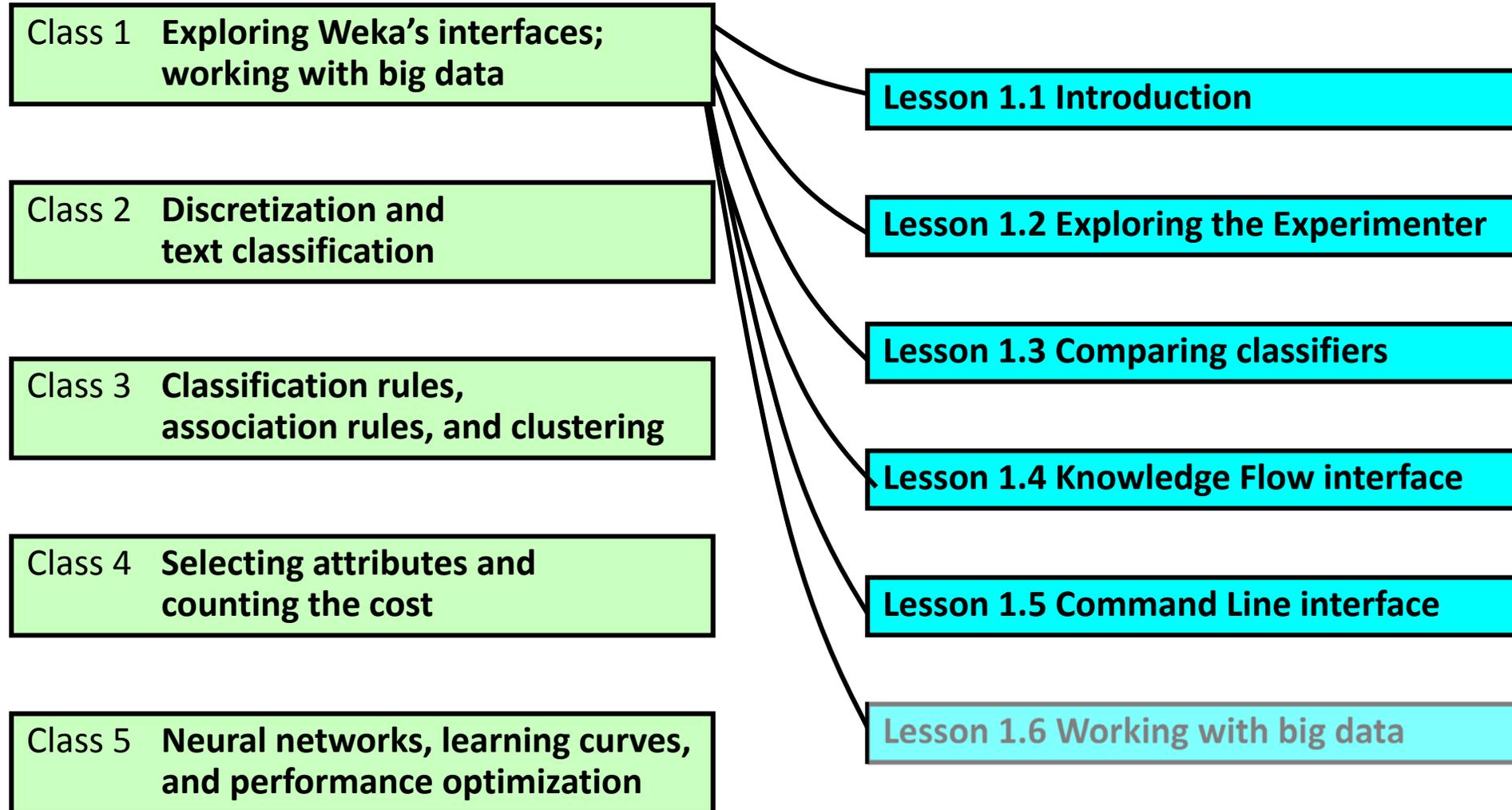
The Command Line interface

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

Lesson 1.5: The Command Line interface



Lesson 1.5: The Command Line interface

Run a classifier from within the CLI

- ❖ Print options for J48:

```
java weka.classifiers.trees.J48
```

- ❖ General options

-h print help info

-t <name of training file> [absolute path name ...]

-T <name of test file>

- ❖ Options specific to J48 (from Explorer configuration panel)

- ❖ Run J48:

```
java weka.classifiers.trees.J48 -C 0.25 -M 2 ←..... copy from Explorer
```

```
-t "C:\Users\ihw\My Documents\Weka datasets\iris.arff" ←..... training set
```

Lesson 1.5: The Command Line interface

Classes and packages

- ❖ J48 is a “class”
 - *a collection of variables, along with some “methods” that operate on them*
- ❖ “Package” is a directory containing related classes

weka.classifiers.trees.J48
←····· ↑····· →·····
packages class

- ❖ Javadoc: the definitive documentation for Weka
Weka-3-6\documentation.html
- ❖ ... find J48 in the “All classes” list

Lesson 1.5: The Command Line interface

Using the Javadoc

“What’s all this geeky stuff?” – *Forget it. Try to ignore things you don’t understand!*

- ❖ Find the “converter” package

weka.core.converters

- ❖ Find the “databaseLoader” class

weka.core.converters.DatabaseLoader

- ❖ Can load from any JDBC database

specify URL, password, SQL query

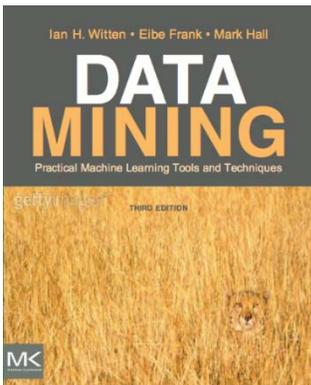
- ❖ It’s in the Explorer’s Preprocess panel, but the documentation is here

Lesson 1.5: The Command Line interface

- ❖ Can do everything the Explorer does from the command line
- ❖ People often open a terminal window instead
 - then you can do scripting (if you know how)
 - ... but you need to set up your environment properly
- ❖ Can copy and paste configured classifiers from the Explorer
- ❖ Advantage: more control over memory usage (next lesson)
- ❖ Javadoc is the definitive source of Weka documentation

Course text

- ❖ Chapter 14 *The Command-Line Interface*





More Data Mining with Weka

Class 1 – Lesson 6

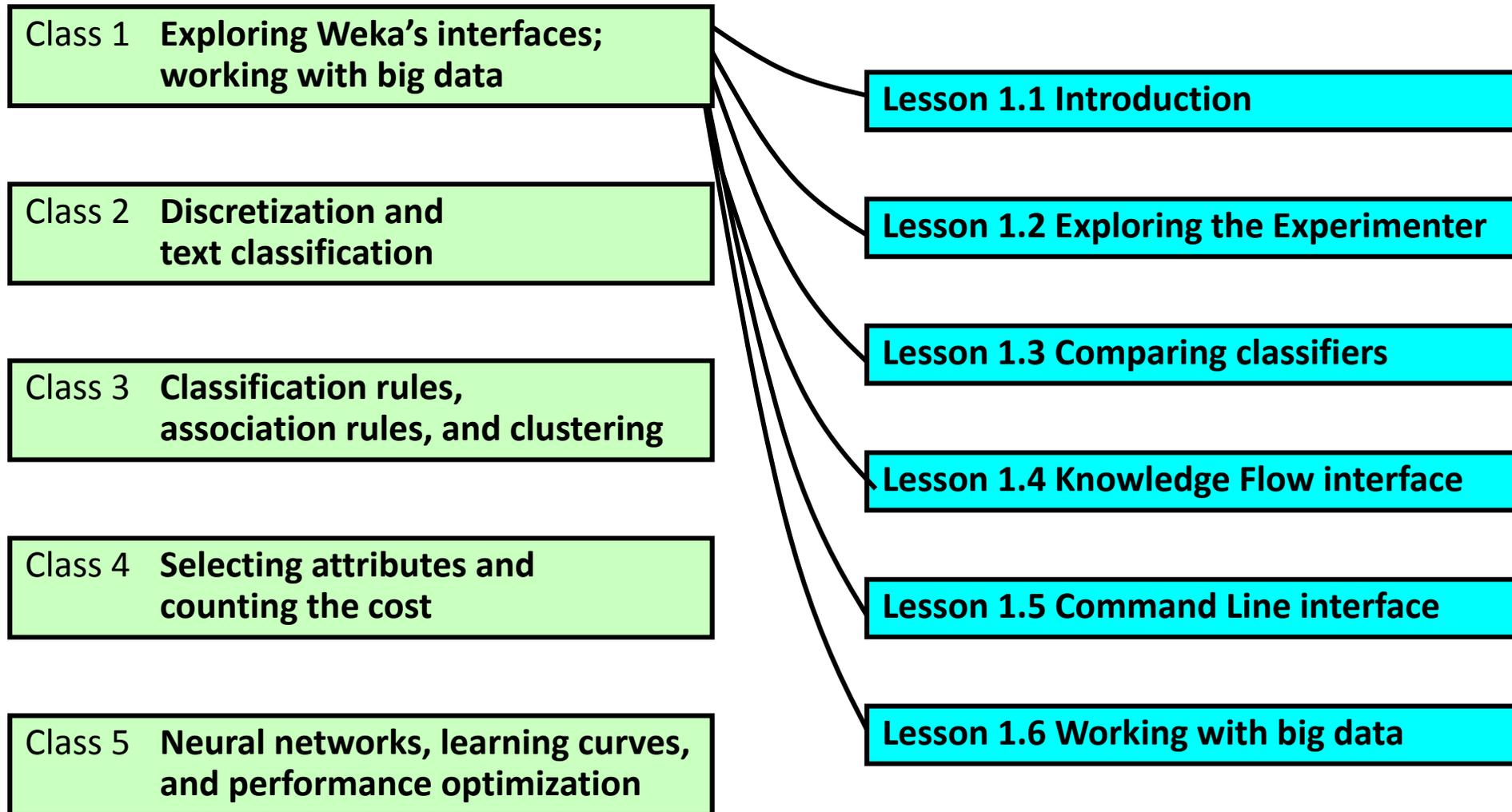
Working with big data

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

Lesson 1.6: Working with big data

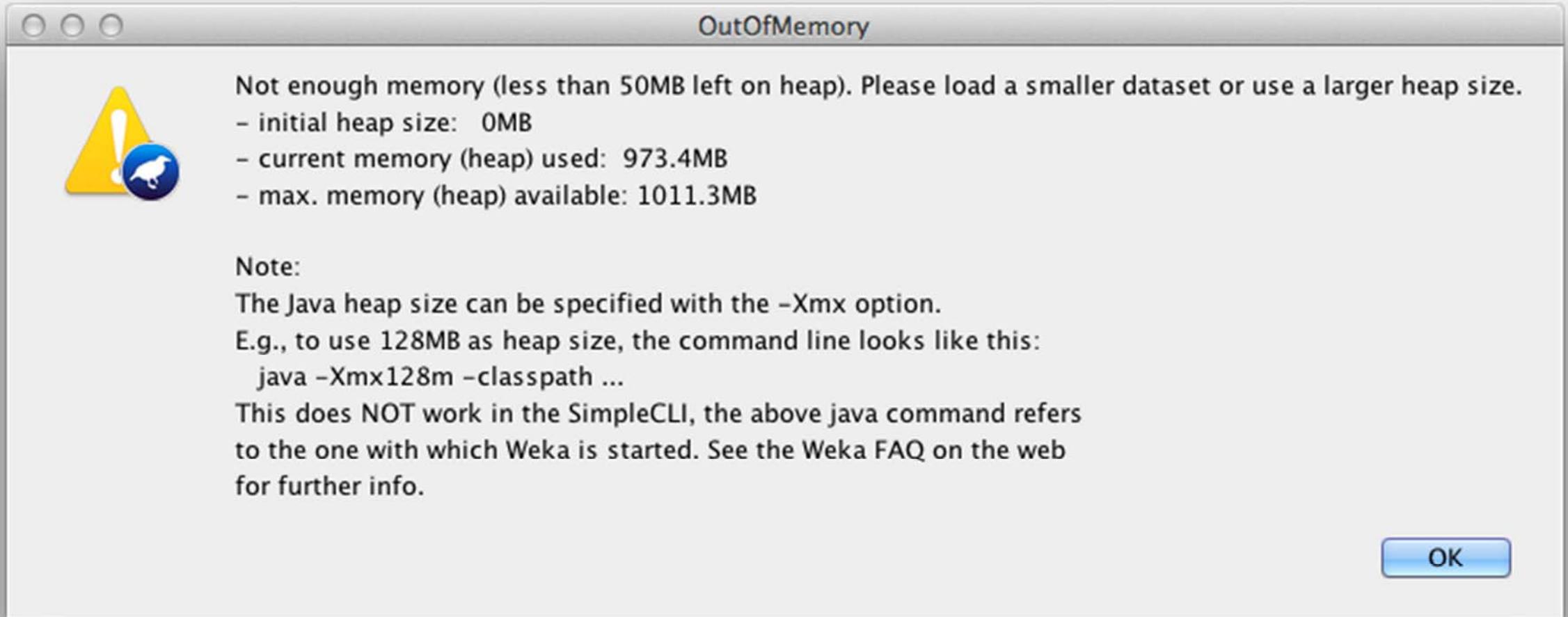


Lesson 1.6: Working with big data

How much can Explorer handle? (~ 1M instances, 25 attributes)

- ❖ Memory information: in Explorer, right-click on “Status”
 - *Free/total/max: 226,366,616 / 236,453,888 / 954,728,448 (bytes) [1 GB]*
 - *Meaning what? Geeks, check out Java’s `freeMemory()`, `totalMemory()`, `maxMemory()` commands*
- ❖ Let’s break it!
- ❖ Download a large dataset?
 - *“covertime” dataset used in the associated Activity*
 - *580,000 instances, 54 attributes (0.75 GB uncompressed)*
- ❖ Weka data generator
 - *Preprocess panel, Generate, choose LED24; show text: 100 instances, 25 attributes*
 - *100,000 examples (use % split!) NaiveBayes 74% J48 73%*
 - *1,000,000 examples NaiveBayes 74% J48 runs out of memory*
 - *2,000,000 examples Generate process grinds to a halt*
- ❖ (Run console version of Weka)

Lesson 1.6: Working with big data



Lesson 1.6: Working with big data

“Updateable” classifiers

- ❖ Incremental classification models: process one instance at a time
 - *AODE, AODEsr, DMNBtext, IB1, IBk, KStar, LWL, NaiveBayesMultinomialUpdateable, NaiveBayesUpdateable, NNge, RacedIncrementalLogitBoost, SPegasos, Winnow*
- ❖ *NaiveBayesUpdateable*: same as *NaiveBayes*
- ❖ *NaiveBayesMultinomialUpdateable*: see lessons on Text Mining
- ❖ *IB1, IBk* (but testing can be very slow)
- ❖ *KStar, LWL* (locally weighted learning): instance-based
- ❖ *SPegasos* (in *functions*)
 - *builds a linear classifier, SVM-style (restricted to numeric or binary class)*
- ❖ *RacedIncrementalLogitBoost*: a kind of boosting

Lesson 1.6: Working with big data

How much can Weka (Simple CLI) handle? – unlimited (conditions apply)

❖ Create a huge dataset

```
java weka.datagenerators.classifiers.classification.LED24 -n 100000 -o C:\Users\ihw\test.arff
```

– Test file with 100 K instances, 5 MB

```
java weka.datagenerators.classifiers.classification.LED24 -n 10000000 -o C:\Users\ihw\train.arff
```

– Training file with 10 M instances; 0.5 GB

❖ Use NaiveBayesUpdateable

```
java weka.classifiers.bayes.NaiveBayesUpdateable -t ...train.arff -T ...test.arff
```

– 74%; 4 mins

– Note: if no test file specified, will do cross-validation, which will fail (non-incremental)

❖ Try with 100 M examples (5 GB training file) – no problem (40 mins)

Lesson 1.6: Working with big data

- ❖ *Explorer* can handle ~ 1M instances, 25 attributes (50 MB file)
- ❖ *Simple CLI* works incrementally wherever it can
- ❖ Some classifier implementations are “Updateable”
 - *find them with Javadoc; see Lesson 1.5 Activity*
- ❖ Updateable classifiers deal with arbitrarily large files (multi GB)
 - *but don't attempt cross-validation*
- ❖ Working with big data can be difficult and frustrating
 - *see the associated Activity*



More Data Mining with Weka

Department of Computer Science
University of Waikato
New Zealand



Creative Commons Attribution 3.0 Unported License



creativecommons.org/licenses/by/3.0/

weka.waikato.ac.nz