

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

Data Mining with Weka

**... a practical course on how to
use Weka for data mining**

**... explains the basic principles
of several popular algorithms**

Ian H. Witten

University of Waikato, New Zealand

Data Mining with Weka

❖ **What's data mining?**

- *We are overwhelmed with data*
- *Data mining is about going from data to information, information that can give you useful predictions*

❖ **Examples??**

- *You're at the supermarket checkout.
You're happy with your bargains ...
... and the supermarket is happy you've bought some more stuff*
- *Say you want a child, but you and your partner can't have one.
Can data mining help?*

❖ **Data mining vs. machine learning**

Data Mining with Weka

❖ **What's Weka?**

– *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



Data Mining with Weka

What will you learn?

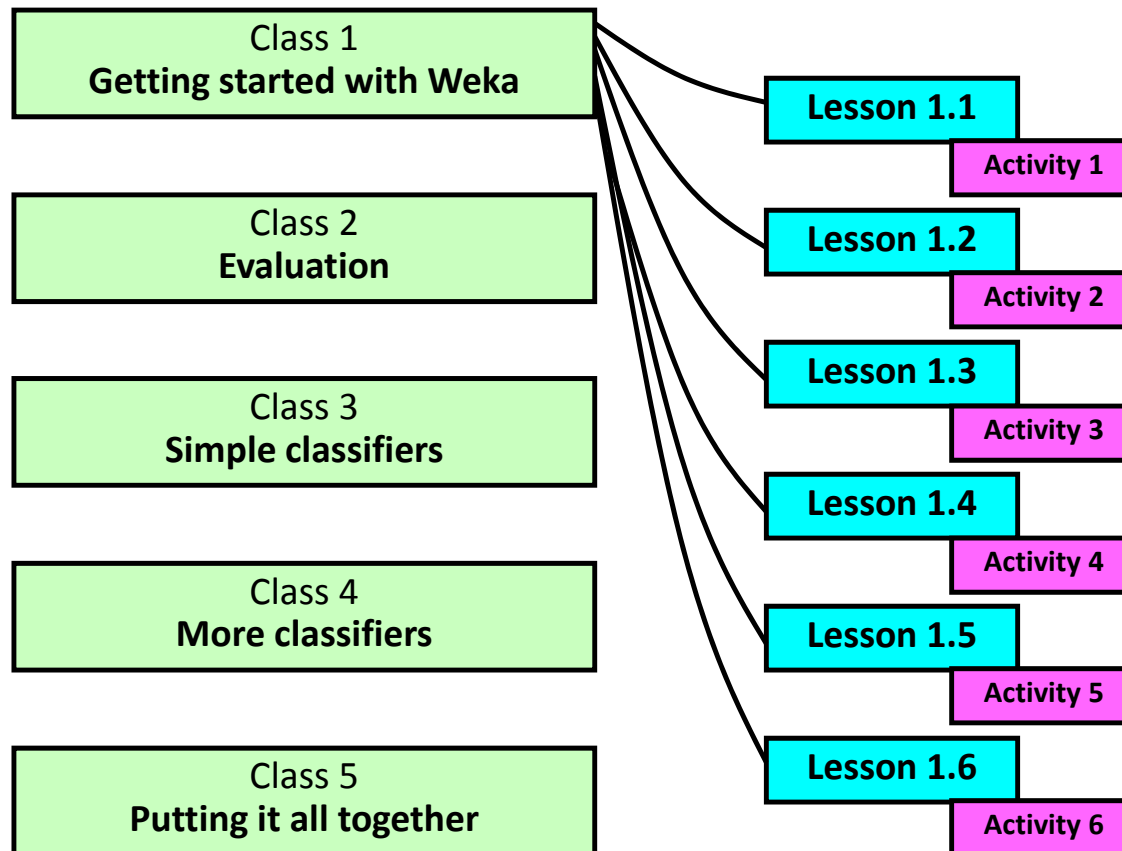
- ❖ Load data into Weka and look at it
- ❖ Use filters to preprocess it
- ❖ Explore it using interactive visualization
- ❖ Apply classification algorithms
- ❖ Interpret the output
- ❖ Understand evaluation methods and their implications
- ❖ Understand various representations for models
- ❖ Explain how popular machine learning algorithms work
- ❖ Be aware of common pitfalls with data mining

**Use Weka on your own data
... and understand what you are doing!**

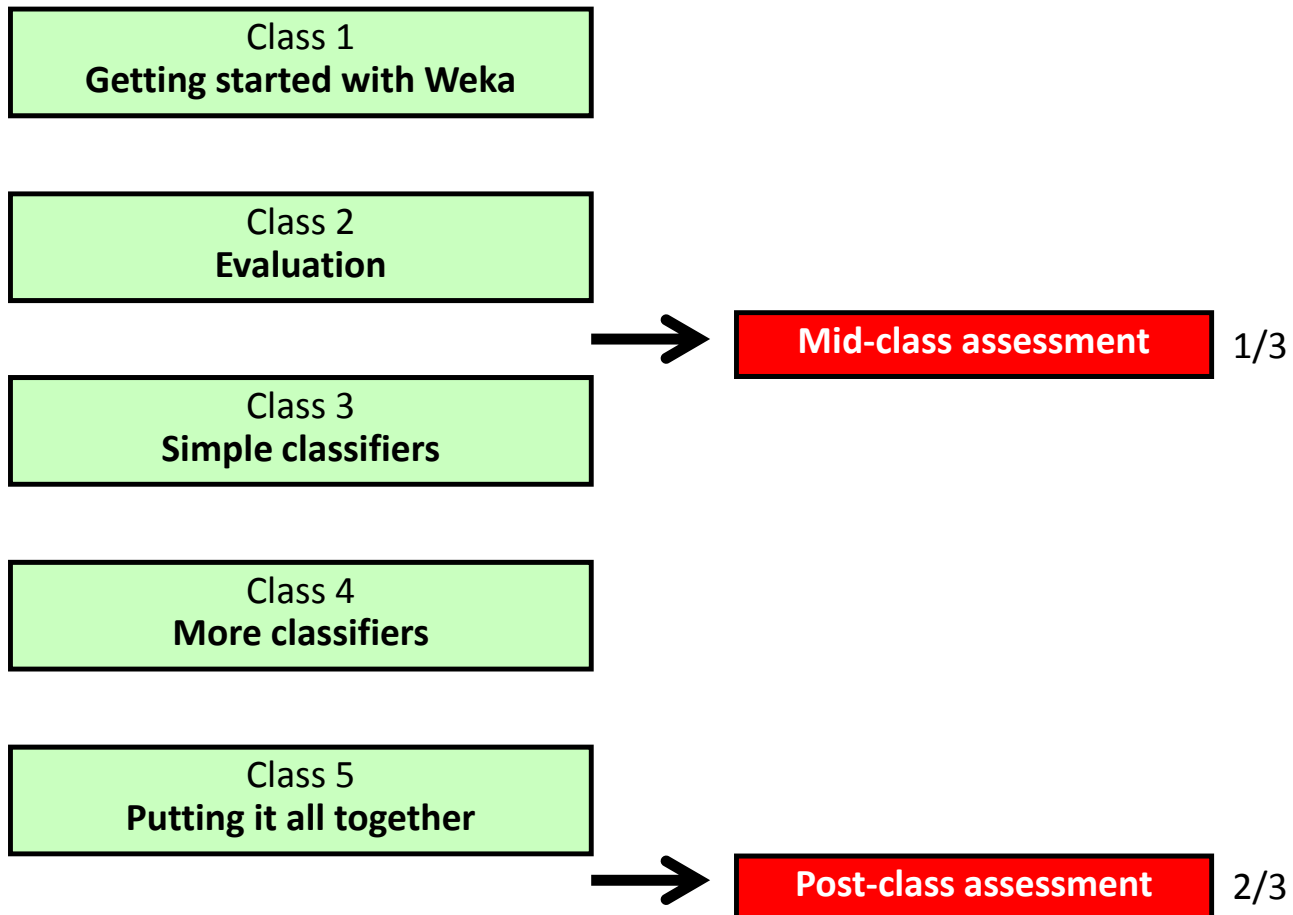
Class 1: Getting started with Weka

- ❖ Install Weka
- ❖ Explore the “Explorer” interface
- ❖ Explore some datasets
- ❖ Build a classifier
- ❖ Interpret the output
- ❖ Use filters
- ❖ Visualize your data set

Course organization



Course organization



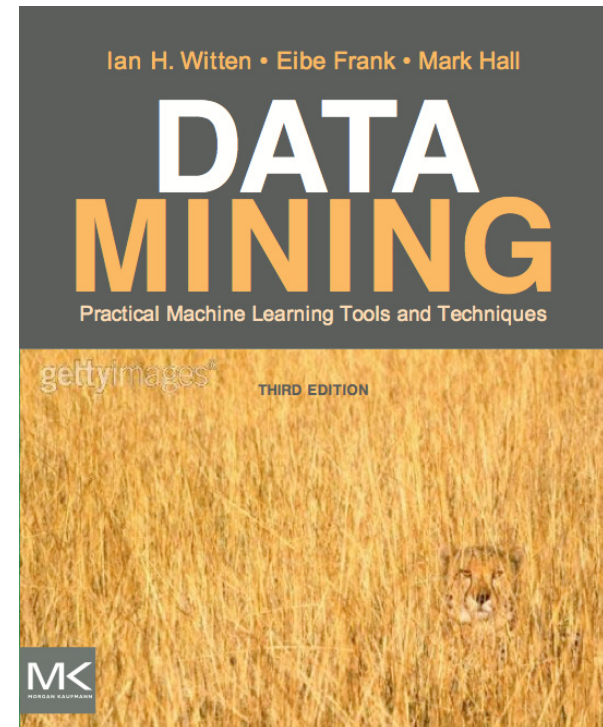
Textbook

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

Data Mining: Practical machine learning tools and techniques,

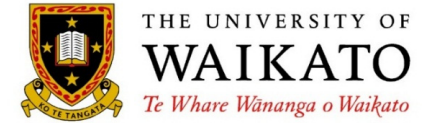
by Ian H. Witten, Eibe Frank and Mark A. Hall. Morgan Kaufmann, 2011

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





World Map by David Niblack, licensed under a Creative Commons Attribution 3.0 Unported License



Data Mining with Weka

Class 1 – Lesson 2

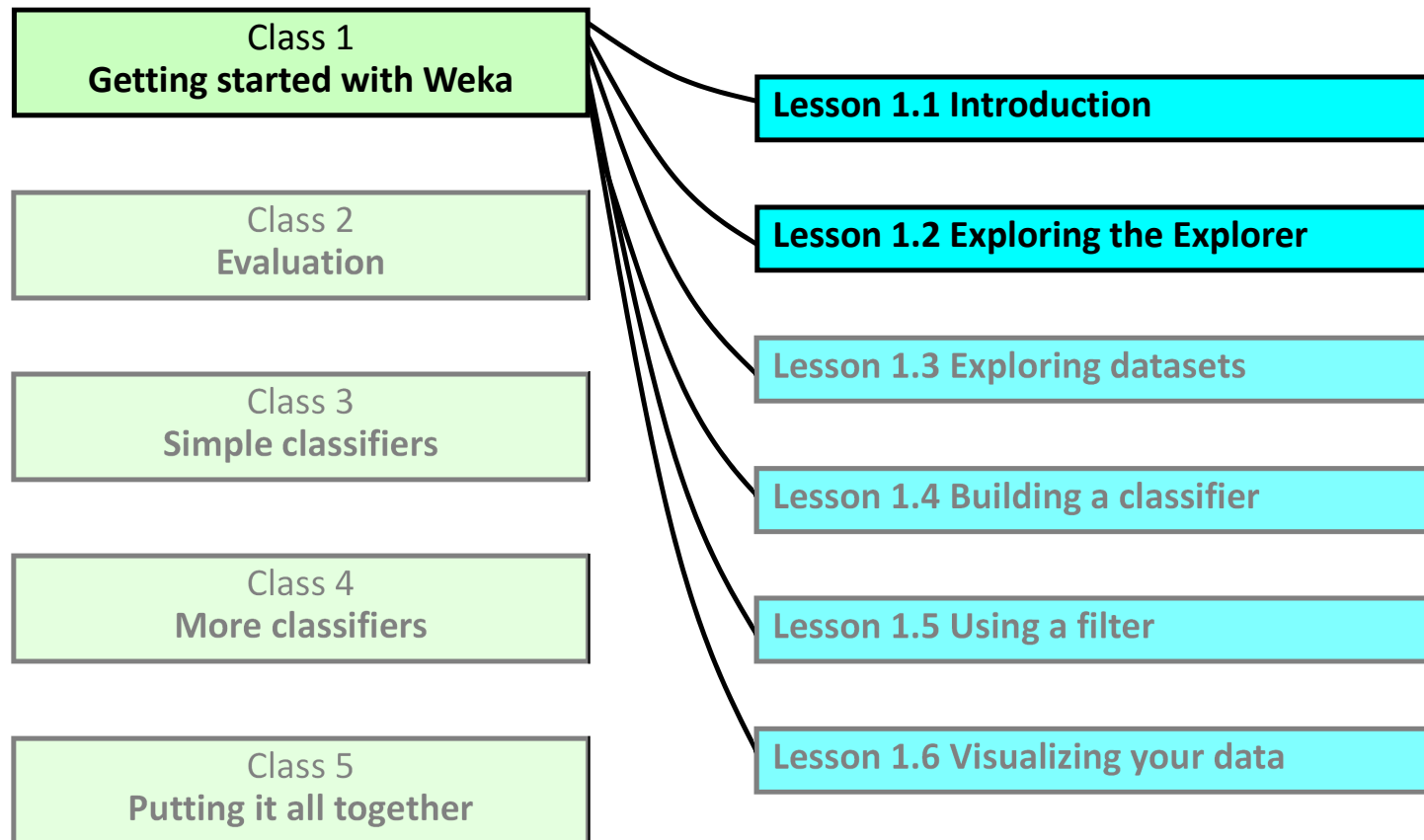
Exploring the Explorer

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

Lesson 1.2: Exploring the Explorer



Lesson 1.2: Exploring the Explorer

Download from

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

(for Windows, Mac, Linux)

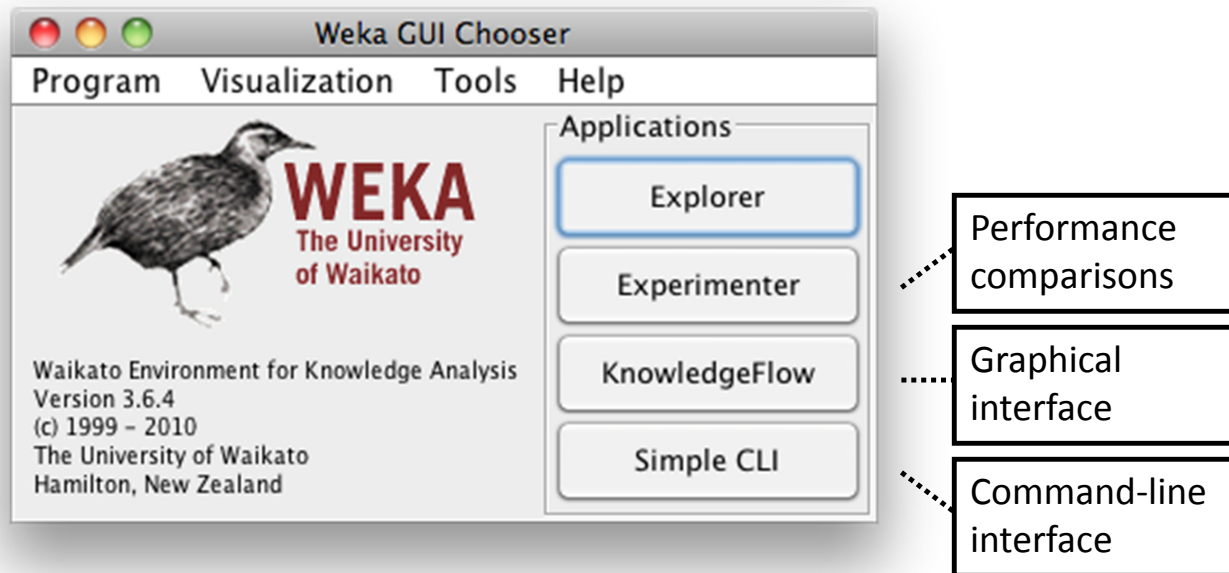
Weka 3.6.10

(the latest stable version of Weka)

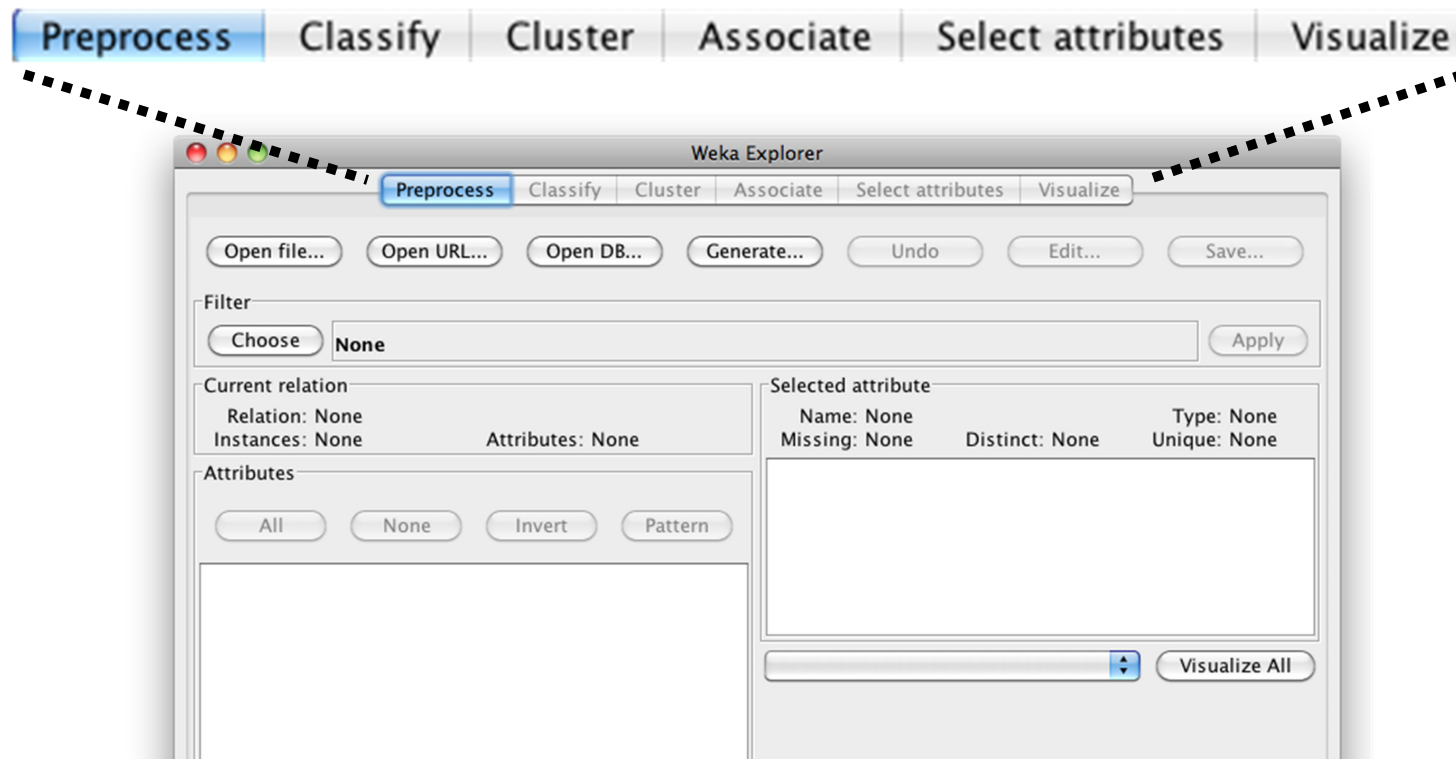
(includes datasets for the course)

(it's important to get the right version, 3.6.10)

Lesson 1.2: Exploring the Explorer



Lesson 1.2: Exploring the Explorer



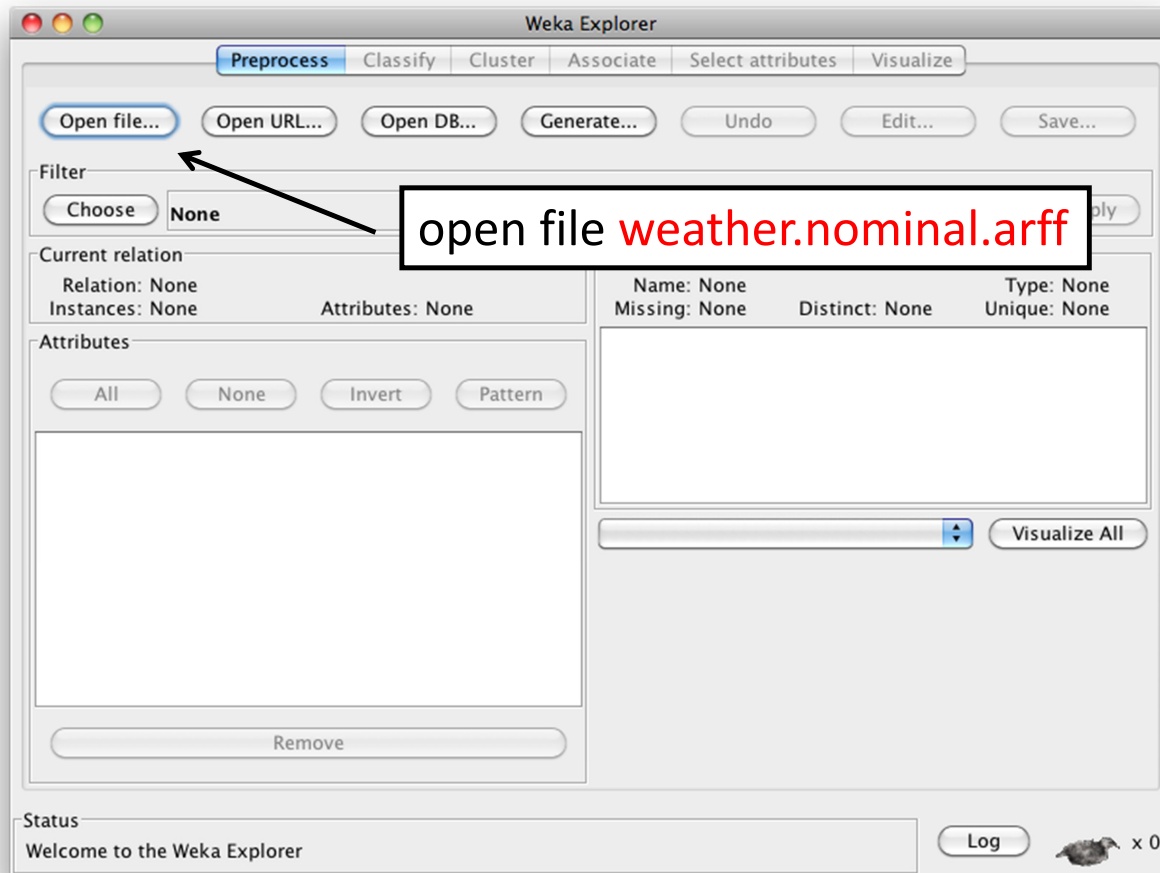
Lesson 1.2: Exploring the Explorer

attributes

instances

	Outlook	Temp	Humidity	Windy	Play
1	Sunny	Hot	High	False	No
2	Sunny	Hot	High	True	No
3	Overcast	Hot	High	False	Yes
4	Rainy	Mild	High	False	Yes
5	Rainy	Cool	Normal	False	Yes
6	Rainy	Cool	Normal	True	No
7	Overcast	Cool	Normal	True	Yes
8	Sunny	Mild	High	False	No
9	Sunny	Cool	Normal	False	Yes
10	Rainy	Mild	Normal	False	Yes
11	Sunny	Mild	Normal	True	Yes
12	Overcast	Mild	High	True	Yes
13	Overcast	Hot	Normal	False	Yes
14	Rainy	Mild	High	True	No

Lesson 1.2: Exploring the Explorer



Lesson 1.2: Exploring the Explorer

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.symbolic
Instances: 14 | Attributes: 5

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

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

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

Class: play (Nom) | Visualize All

Status: OK | Log | x 0

attributes

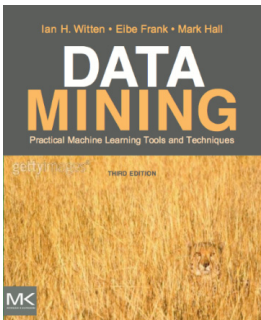
attribute values

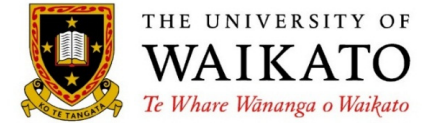
Lesson 1.2: Exploring the Explorer

- ❖ Install Weka
- ❖ Get datasets
- ❖ Open Explorer
- ❖ Open a dataset (*weather.nominal.arff*)
- ❖ Look at attributes and their values
- ❖ Edit the dataset
- ❖ Save it?

Course text

- ❖ Section 1.2 *The weather problem*
- ❖ Chapter 10 *Introduction to Weka*





Data Mining with Weka

Class 1 – Lesson 3

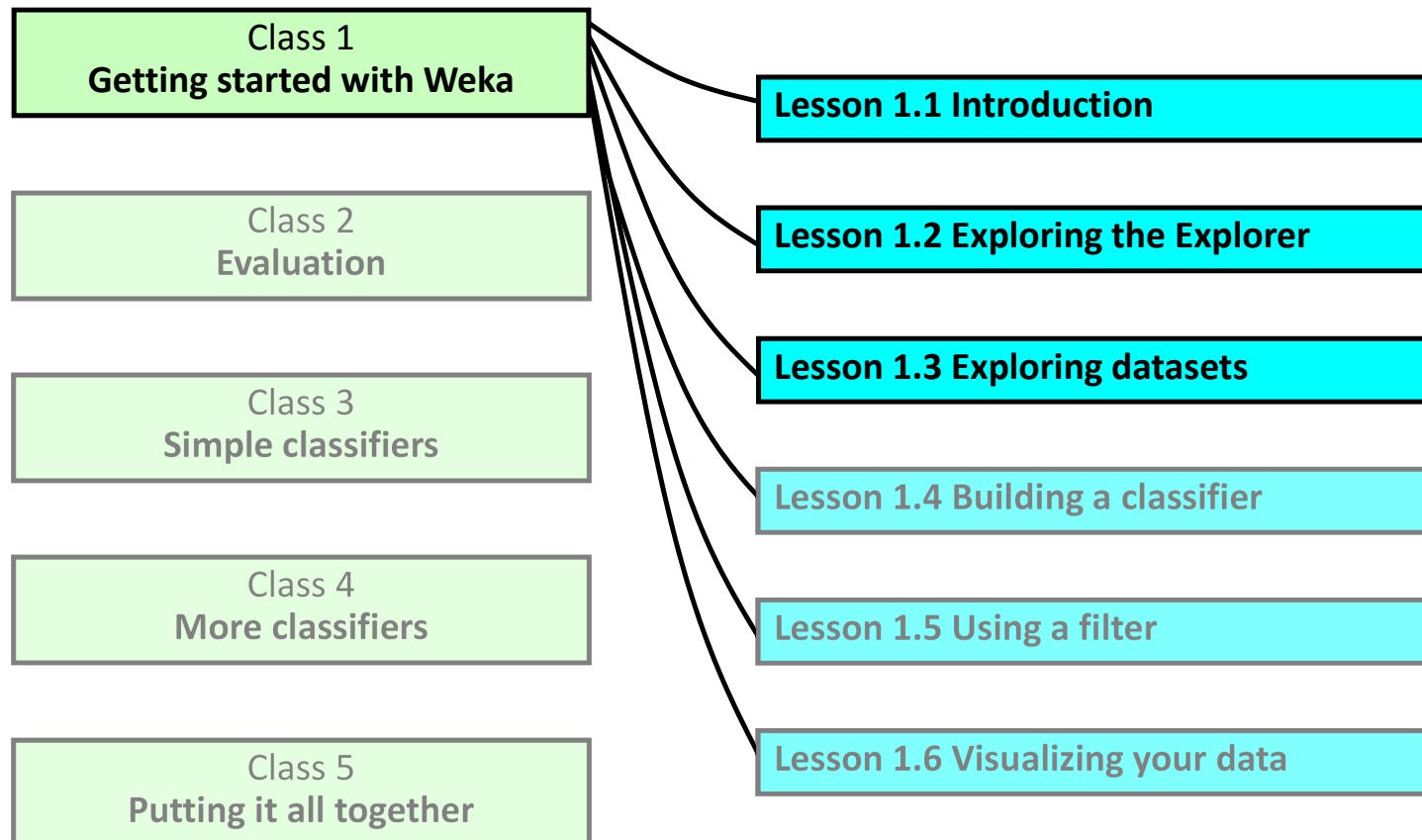
Exploring datasets

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

Lesson 1.3: Exploring datasets



Lesson 1.3: Exploring datasets

		attributes				
		Outlook	Temp	Humidity	Windy	Play
instances	1	Sunny	Hot	High	False	No
	2	Sunny	Hot	High	True	No
	3	Overcast			False	Yes
	4	Rainy			False	Yes
	5	Partly cloudy		Normal	False	Yes
	6	Partly cloudy		Normal	True	No
	7	Partly cloudy	Cool	Normal	True	Yes
	8	Partly cloudy	Mild			No
	9	Sunny	Cool	Normal	False	Yes
	10	Rainy	Mild	Normal	False	Yes
	11	Sunny	Mild	Normal	True	Yes
	12	Overcast	Mild	High	True	Yes
	13	Overcast	Hot	Normal	False	Yes
	14	Rainy	Mild	High	True	No

**Classification problem:
predict the "class" value**

Lesson 1.3: Exploring datasets

The screenshot shows the Weka Explorer interface with the 'Preprocess' tab selected. The 'Open file...' button is highlighted with a callout box containing the text 'open file weather.nominal.arff'. The 'Attributes' list on the left includes 'outlook', 'temperature', 'humidity', 'windy', and 'play', with 'outlook' selected. A callout box labeled 'attributes' points to this list, and another callout box labeled 'class' points to the 'play' attribute. The 'Current relation' section shows 'Relation: weather.symbolic' and 'Instances: 14'. The 'Attributes' section shows 'Attributes: 5'. The 'outlook' attribute details are shown in a table:

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

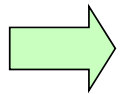
A callout box labeled 'attribute values' points to this table. Below the table, the 'Class: play (Nom)' is shown with a dropdown menu and a 'Visualize All' button. The visualization area shows three stacked bar charts for the 'outlook' attribute, with the first bar having a total height of 5, the second 4, and the third 5. The status bar at the bottom shows 'Status OK' and a 'Log' button.

Lesson 1.3: Exploring datasets

Classification

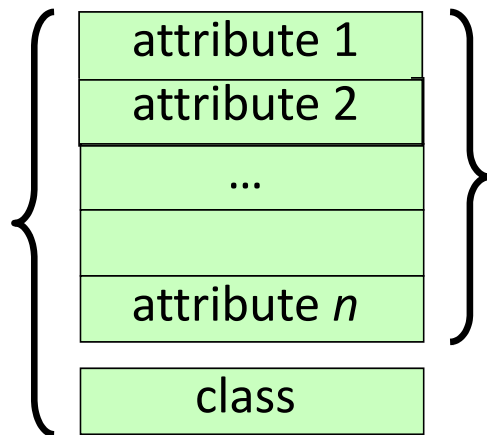
sometimes called “supervised learning”

Dataset: classified examples



“Model” that classifies new examples

classified example



instance:
fixed set of features

discrete (“nominal”)
continuous (“numeric”)

discrete: “classification” problem
continuous: “regression” problem

Lesson 1.3: Exploring datasets

The screenshot shows the Weka Explorer application window. The 'Preprocess' tab is active. The 'Open file...' button is highlighted with a callout box containing the text 'open file weather.numeric.arff'. The 'Attributes' section on the left lists five attributes: outlook, temperature, humidity, windy, and play. The 'play' attribute is selected, and a callout box labeled 'class' points to it. The 'Current relation' section shows 'Relation: weather.symbolic' and 'Instances: 14'. The 'Attributes: 5' section shows a table with three columns: No., Label, and Count. The table contains the following data:

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

A callout box labeled 'attribute values' points to this table. Below the table, the 'Class: play (Nom)' is selected, and a 'Visualize All' button is visible. The visualization area shows three stacked bar charts, each representing a different outlook value (sunny, overcast, rainy). Each bar is divided into two segments: a blue segment at the bottom and a red segment at the top. The total height of each bar corresponds to the count of instances for that outlook value (5 for sunny, 4 for overcast, and 5 for rainy). The status bar at the bottom shows 'Status OK' and a 'Log' button.

Lesson 1.3: Exploring datasets

The screenshot shows the Weka Explorer application window. The 'Preprocess' tab is active. The 'Open file...' button is highlighted with a red box and an arrow pointing to it, with the text 'open file glass.arff' written inside the box. The 'Filter' section shows 'None' selected. The 'Current relation' section shows 'Relation: Glass' and 'Instances: 214'. The 'Attributes' section shows a list of 10 attributes: RI, Na, Mg, Al, Si, K, Ca, Ba, Fe, and Type. The 'RI' attribute is selected. The 'Statistic' table shows the following values:

Statistic	Value
Minimum	1.511
Maximum	1.534
Mean	1.518
StdDev	0.003

The 'Class: Type (Nom)' dropdown is set to 'Type (Nom)'. A histogram is displayed below the dropdown, showing the distribution of the 'Type' attribute. The histogram has 10 bins, with the following counts: 3, 4, 39, 84, 39, 16, 17, 4, 3, 3, 0, 1, 1. The x-axis ranges from 1.51 to 1.53.

Status: OK

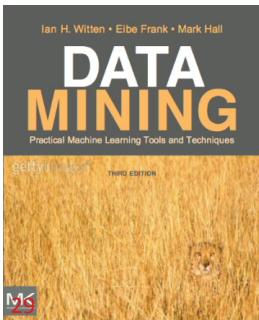
Log x 0

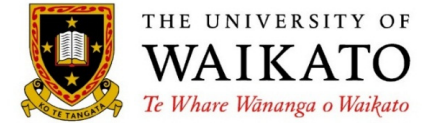
Lesson 1.3: Exploring datasets

- ❖ The classification problem
- ❖ *weather.nominal, weather.numeric*
- ❖ Nominal vs numeric attributes
- ❖ ARFF file format
- ❖ *glass.arff* dataset
- ❖ Sanity checking attributes

Course text

- ❖ Section 11.1 *Preparing the data*
Loading the data into the Explorer





Data Mining with Weka

Class 1 – Lesson 4

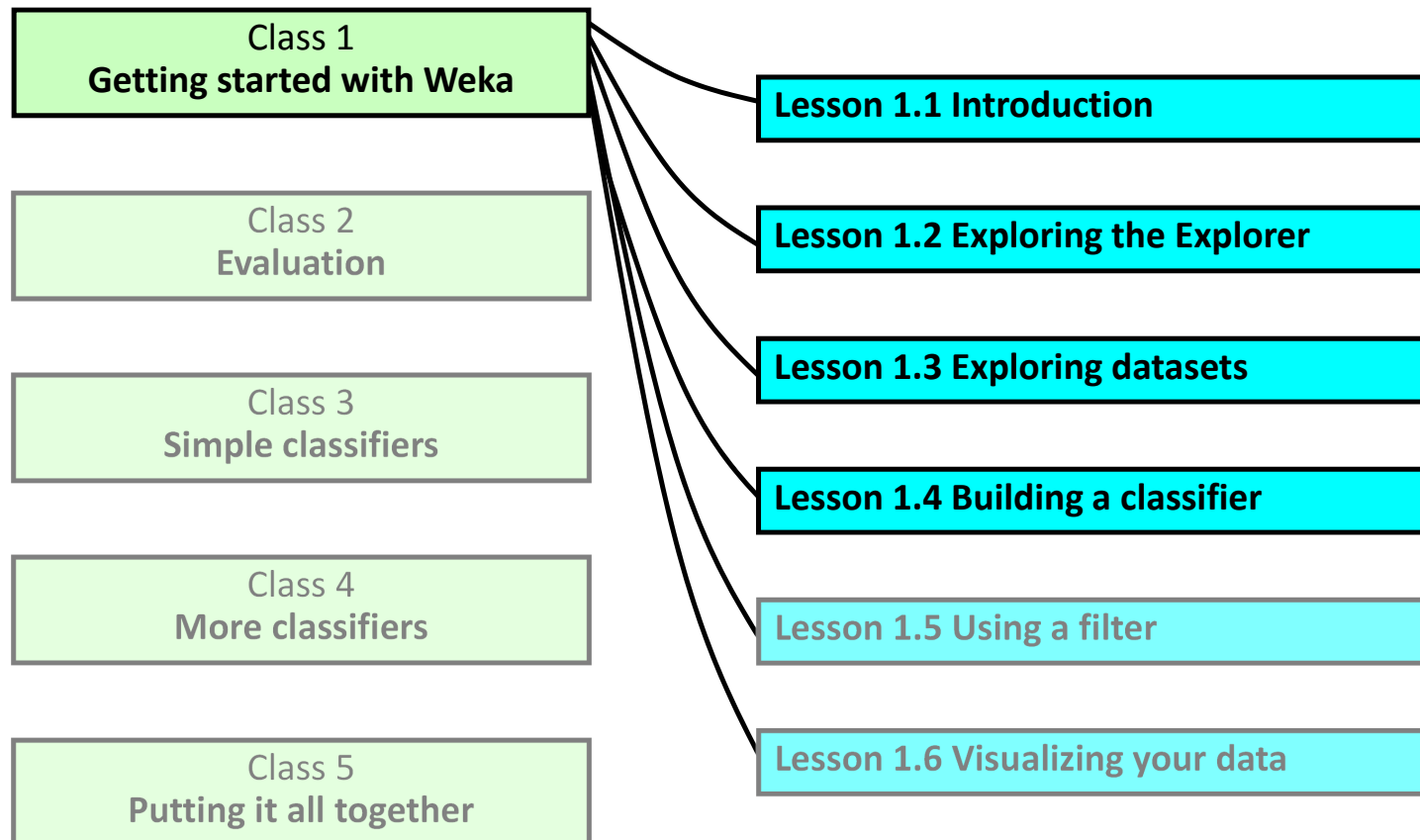
Building a classifier

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

Lesson 1.4: Building a classifier



Lesson 1.4: Building a classifier

Use J48 to analyze the glass dataset

- ❖ Open file `glass.arff`
(or leave it open from the last lesson)
- ❖ Check the available classifiers
- ❖ Choose the J48 decision tree learner (`trees>J48`)
- ❖ Run it
- ❖ Examine the output
- ❖ Look at the correctly classified instances
... and the confusion matrix

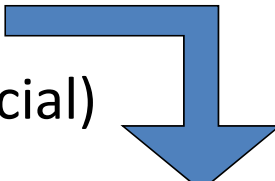
Lesson 1.4: Building a classifier

Investigate J48

- ❖ Open the configuration panel
- ❖ Check the *More* information
- ❖ Examine the options
- ❖ Use an unpruned tree
- ❖ Look at leaf sizes
- ❖ Set **minNumObj** to 15 to avoid small leaves
- ❖ Visualize tree using right-click menu

Lesson 1.4: Building a classifier

From C4.5 to J48

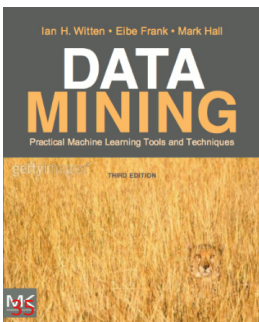
- ❖ ID3 (1979)
 - ❖ **C4.5** (1993)
 - ❖ C4.8 (1996?)
 - ❖ C5.0 (commercial)
- 
- J48**

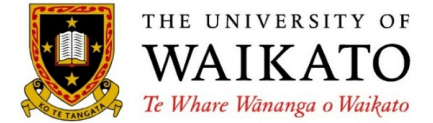
Lesson 1.4: Building a classifier

- ❖ Classifiers in Weka
- ❖ Classifying the *glass* dataset
- ❖ Interpreting J48 output
- ❖ J48 configuration panel
- ❖ ... option: pruned vs unpruned trees
- ❖ ... option: avoid small leaves
- ❖ J48 ~ C4.5

Course text

- ❖ Section 11.1 *Building a decision tree*
Examining the output





Data Mining with Weka

Class 1 – Lesson 5

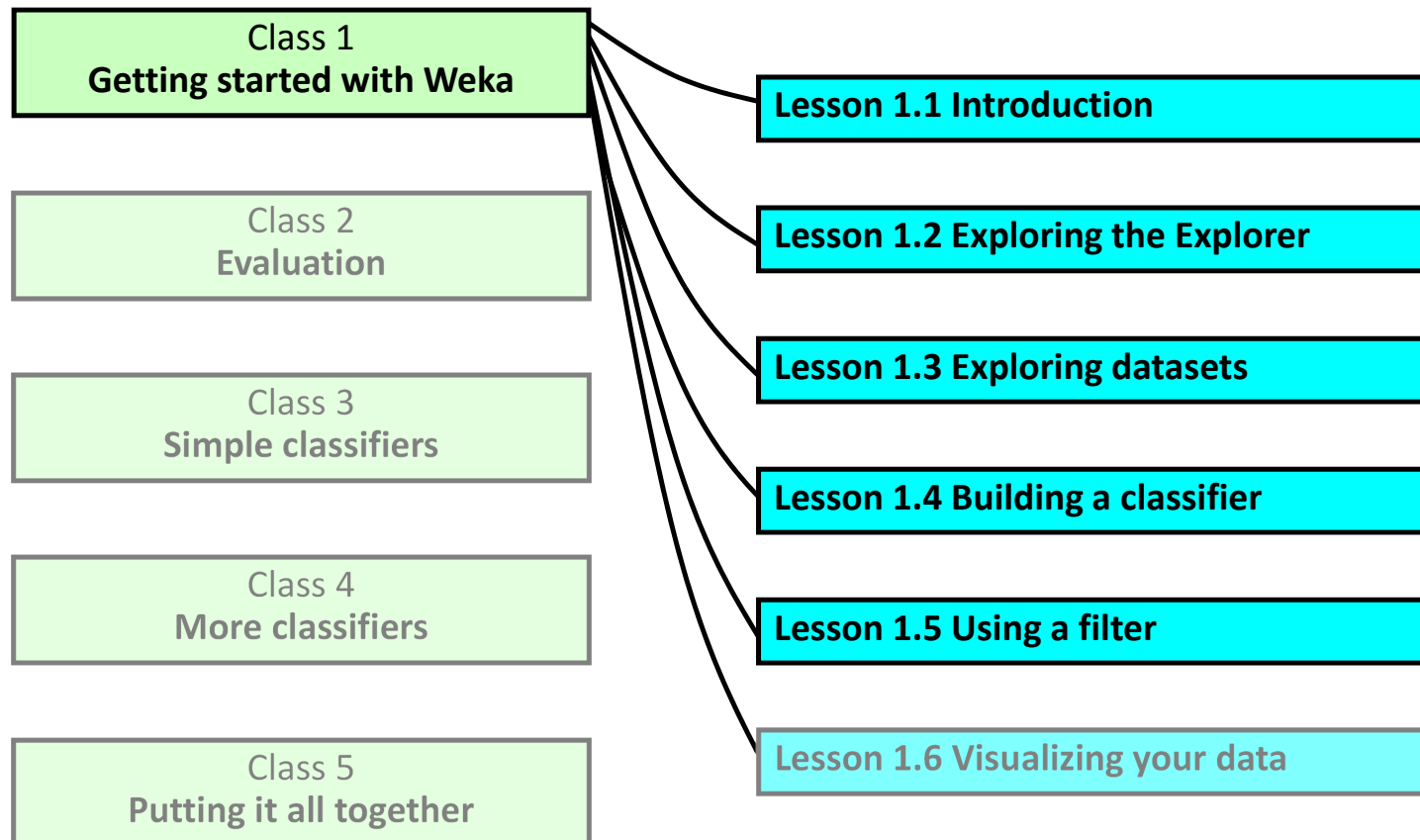
Using a filter

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

Lesson 1.5: Using a filter



Lesson 1.5: Using a filter

Use a filter to remove an attribute

- ❖ Open **weather.nominal.arff** (again!)
- ❖ Check the filters
 - supervised vs unsupervised
 - attribute vs instance
- ❖ Choose the **unsupervised attribute** filter *Remove*
- ❖ Check the *More* information; look at the options
- ❖ Set **attributeIndices** to **3** and click OK
- ❖ Apply the filter
- ❖ Recall that you can *Save* the result
- ❖ Press *Undo*

Lesson 1.5: Using a filter

Remove instances where *humidity* is *high*

- ❖ Supervised or unsupervised?
- ❖ Attribute or instance?
- ❖ Look at them
- ❖ Select *RemoveWithValues*
- ❖ Set *attributeIndex*
- ❖ Set *nominalIndices*
- ❖ Apply
- ❖ *Undo*

Lesson 1.5: Using a filter

Fewer attributes, better classification!

- ❖ Open `glass.arff`
- ❖ Run J48 (`trees>J48`)
- ❖ Remove Fe
- ❖ Remove all attributes except RI and MG
- ❖ Look at the decision trees

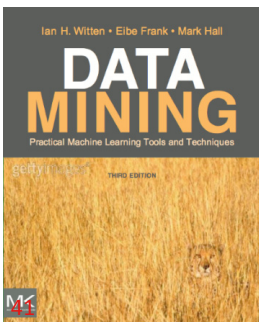
- ❖ Use right-click menu to visualize decision trees

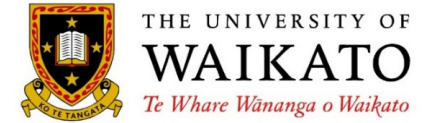
Lesson 1.5: Using a filter

- ❖ Filters in Weka
- ❖ Supervised vs unsupervised, attribute vs instance
- ❖ To find the right one, you need to look!
- ❖ Filters can be very powerful
- ❖ Judiciously removing attributes can
 - improve performance
 - increase comprehensibility

Course text

- ❖ Section 11.2 *Loading and filtering files*





Data Mining with Weka

Class 1 – Lesson 6

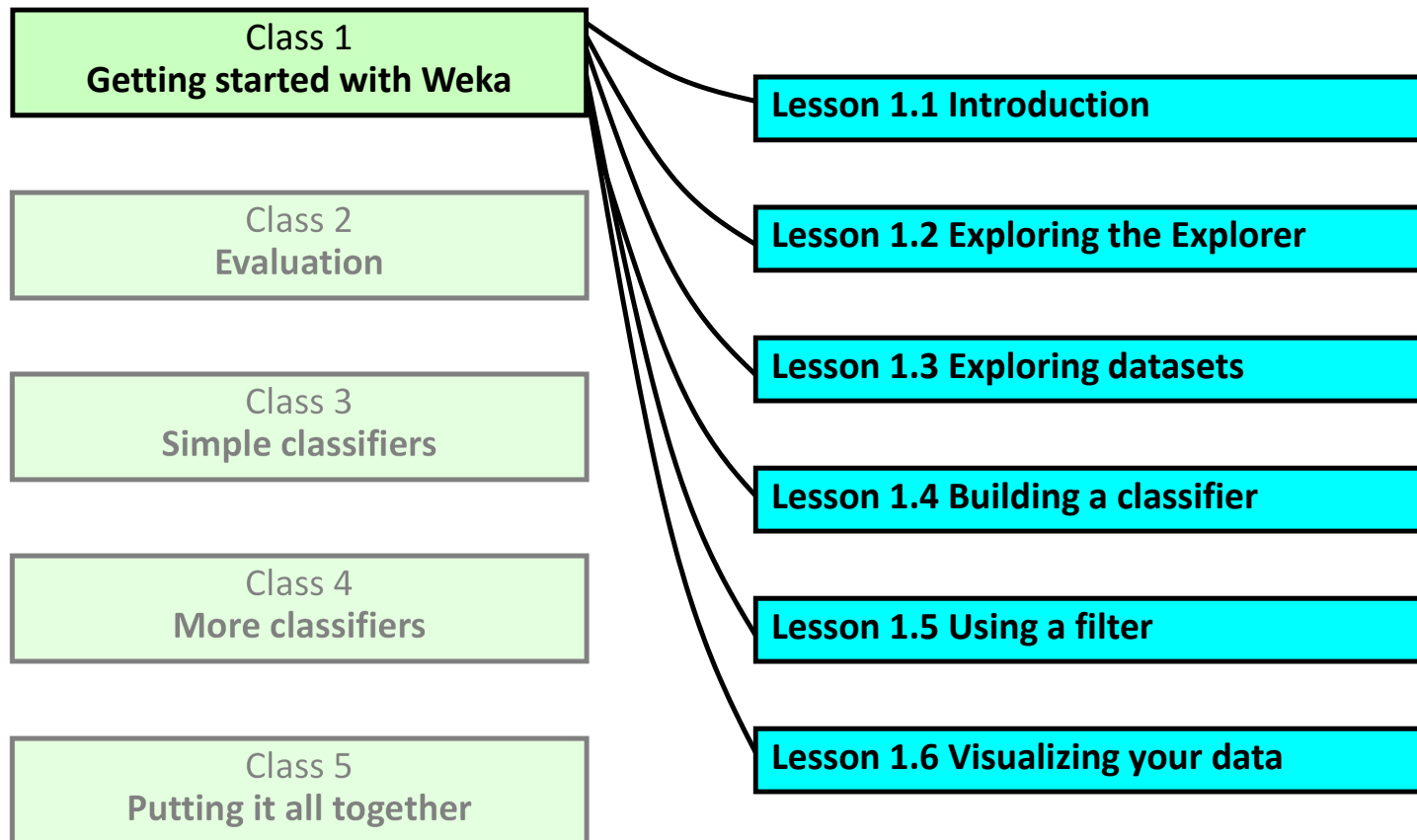
Visualizing your data

Ian H. Witten

Department of Computer Science
University of Waikato
New Zealand

weka.waikato.ac.nz

Lesson 1.6: Visualizing your data



Lesson 1.6: Visualizing your data

Using the Visualize panel

- ❖ Open iris.arff
- ❖ Bring up Visualize panel
- ❖ Click one of the plots; examine some instances
- ❖ Set x axis to petalwidth and y axis to petallength
- ❖ Click on Class colour to change the colour
- ❖ Bars on the right change correspond to attributes: click for x axis; right-click for y axis
- ❖ Jitter slider
- ❖ Show Select Instance: Rectangle option
- ❖ Submit, Reset, Clear and Save

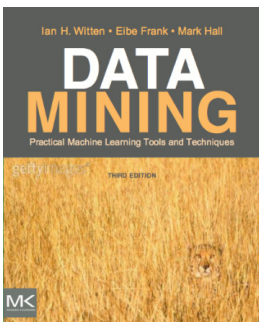
Lesson 1.6: Visualizing your data

Visualizing classification errors

- ❖ Run J48 (trees>J48)
- ❖ Visualize classifier errors (from Results list)
- ❖ Plot predictedclass against class
- ❖ Identify errors shown by confusion matrix

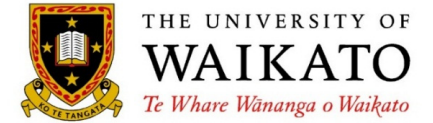
Lesson 1.6: Visualizing your data

- ❖ Get down and dirty with your data
- ❖ Visualize it
- ❖ Clean it up by deleting outliers
- ❖ Look at classification errors
 - (there's a filter that allows you to add classifications as a new attribute)



Course text

Section 11.2 Visualization



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