



Weka



UNIVERSITY OF
BIRMINGHAM

Introduction to WEKA

Lecturer: Dr. Sharu Theresa Jose



Learning Outcomes

- Familiarize with a popular ML software – WEKA
- Familiarize with different data sets
- Application of ML algorithms to real-world data sets



WEKA

- WEKA = Waikato Environment for Knowledge Analysis
- Developed at the University of Waikato in New Zealand
- Java-based open-source software
- Collection of machine learning algorithms and data pre-processing tools



What's in WEKA?

- Provides implementations of learning algorithms that can be easily applied on data set (without actually coding !)
- Includes methods for main ML problems – classification, regression, clustering, attribute selection
- How can we use it?
 - Apply a learning method to a dataset and analyze its output
 - Learn models to predict on new instances
 - Compare performances of different learning methods



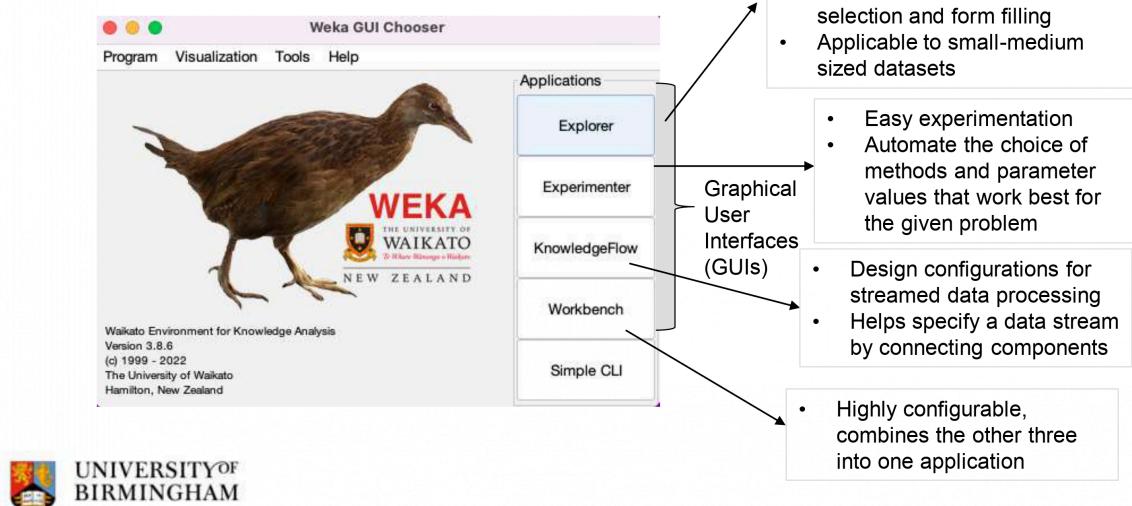
Downloading and using Weka

- Can be downloaded from <https://www.cs.waikato.ac.nz/ml/weka/>
- Access Weka via school's jupyter notebook server <https://jupyterhub.oc1.aws.cs.bham.ac.uk>. (This, however, does not support graphical user interface. Details on how to work with Jupyter notebook server will be uploaded as a separate video.)



Starting WEKA

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



The Explorer

- Preprocess**: choose dataset and modify it in various ways
- Classify**: train learning schemes that perform classification or regression and evaluate them
- Cluster**: learn clusters from the dataset
- Associate**: learn association rules for the data and evaluate them
- Select Attributes**: select the most relevant aspects of the dataset
- Visualize**: view different two-dimensional plots of the data and interact with them



Getting the Explorer Started: Data Format

- Data storage format: ARFF (Attribute-Relation File Format)

A	B	C	D	E
1 outlook	temperature	humidity	windy	play
2 sunny	85	85	FALSE	no
3 sunny	80	90	TRUE	no
4 overcast	83	86	FALSE	yes
5 rainy	70	96	FALSE	yes
6 rainy	68	80	FALSE	yes
7 rainy	65	70	TRUE	no
8 overcast	64	65	TRUE	yes
9 sunny	72	95	FALSE	no
10 sunny	69	70	FALSE	yes
11 rainy	75	80	FALSE	yes
12 sunny	75	70	TRUE	yes
13 overcast	72	90	TRUE	yes
14 overcast	81	75	FALSE	yes
15 rainy	71	91	TRUE	no

(a) Spreadsheet.

```
outlook,temperature,humidity,windy,play
sunny,85,85, FALSE,no
sunny,80,90, TRUE,no
overcast,83,86, FALSE, yes
rainy,70,96, FALSE, yes
rainy,68,80, FALSE, yes
rainy,65,70, TRUE, no
overcast,64,65, TRUE, yes
sunny,72,95, FALSE, no
sunny,69,70, FALSE, yes
rainy,75,80, FALSE, yes
sunny,75,70, TRUE, yes
overcast,72,90, TRUE, yes
overcast,81,75, FALSE, yes
rainy,71,91, TRUE, no
```

(b) CSV.

Name of data set
Attribute Information
Data instances

```
@relation weather
@attribute outlook {sunny, overcast, rainy}
@attribute humidity numeric
@attribute windy {TRUE, FALSE}
@attribute play {yes, no}

@data
sunny,85,85, FALSE,no
sunny,80,90, TRUE,no
overcast,83,86, FALSE, yes
rainy,70,96, FALSE, yes
rainy,68,80, FALSE, yes
rainy,65,70, TRUE, no
overcast,64,65, TRUE, yes
sunny,72,95, FALSE, no
sunny,69,70, FALSE, yes
rainy,75,80, FALSE, yes
sunny,75,70, TRUE, yes
overcast,72,90, TRUE, yes
overcast,81,75, FALSE, yes
rainy,71,91, TRUE, no
```

(c) ARFF.



UNIVERSITY OF
BIRMINGHAM

Demo



UNIVERSITY OF
BIRMINGHAM

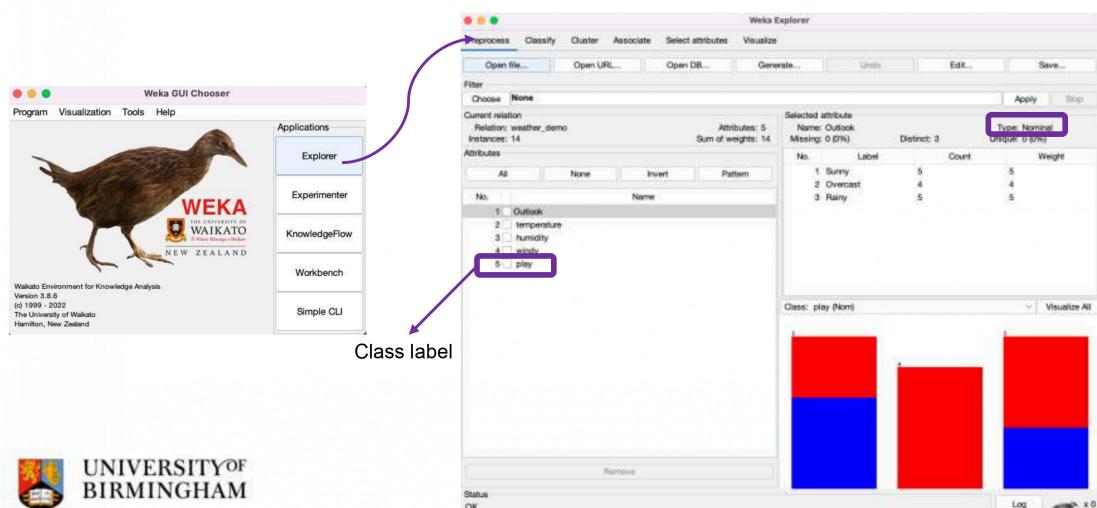
Explorer → Preprocess

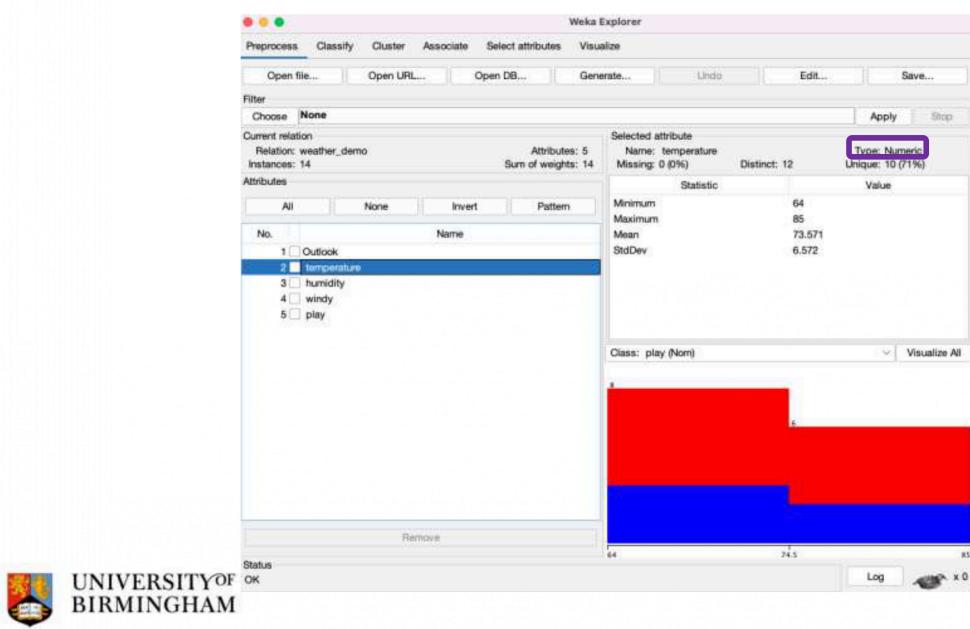
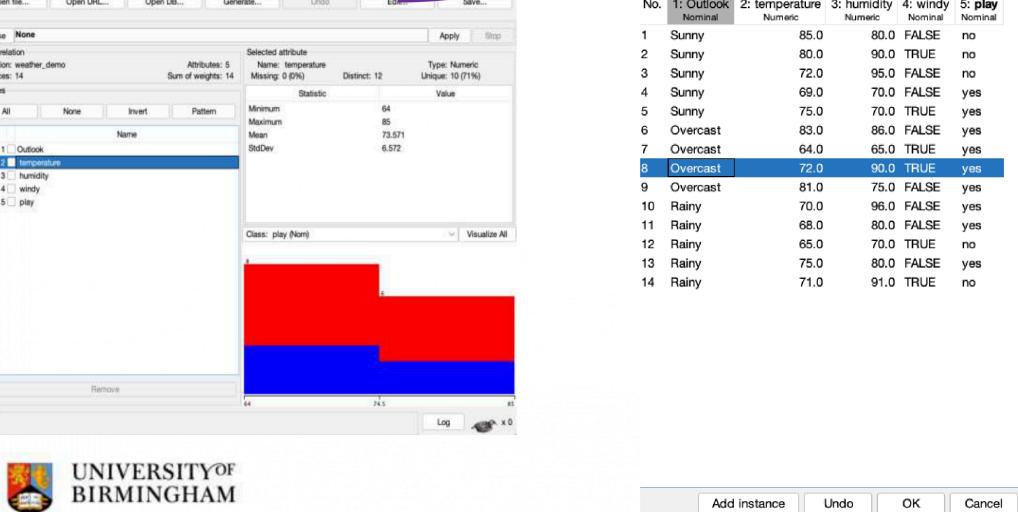
- Load data files
- Pre-processing tools or ‘filters’



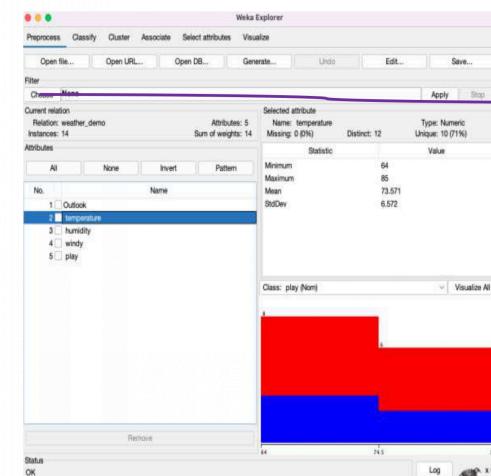
Loading data

- Preprocess → Open file → select .arff file from the desired directory

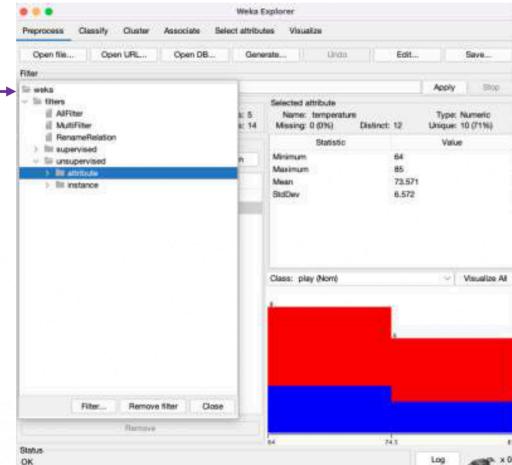




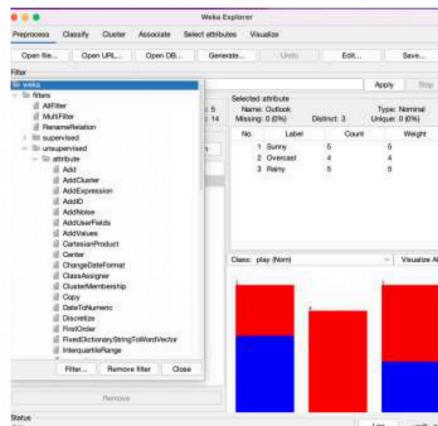
Pre-processing via Filters



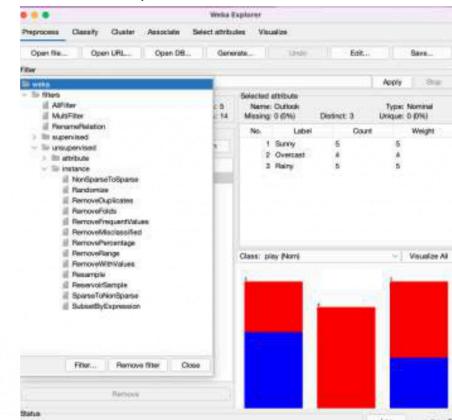
- Supervised filters – consider ‘class’ value
- Unsupervised filters – ignore ‘class’ value



- Unsupervised attribute filters – affect specific attributes or a set of attributes
 - Example: Discretize, DataToNumeric, StringToWordVector



- Unsupervised instance filters – affect all instances in a data set
 - Randomize : randomizes the order of instances in the dataset
 - Resample : produce a random sample by sampling with or without replacement



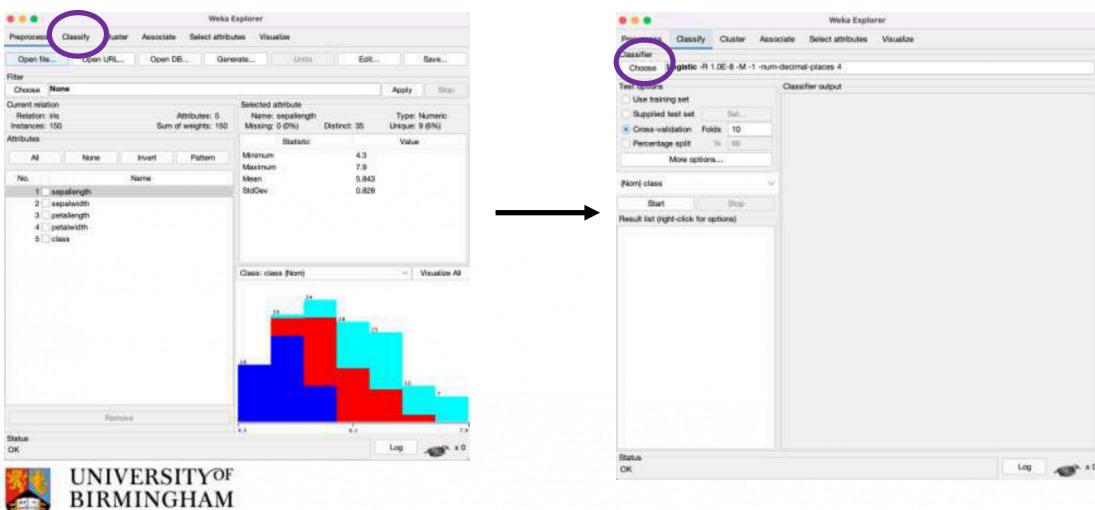
Explorer → Classify

Includes various classification algorithms:

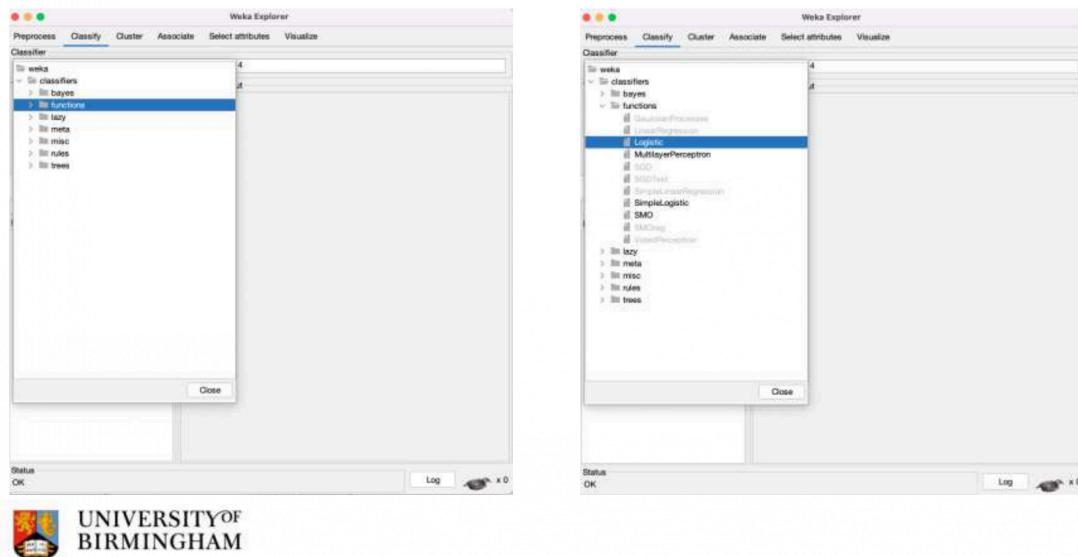
1. Logistic Regression
2. K-Nearest Neighbors
3. Decision Trees
4. Naïve Bayes



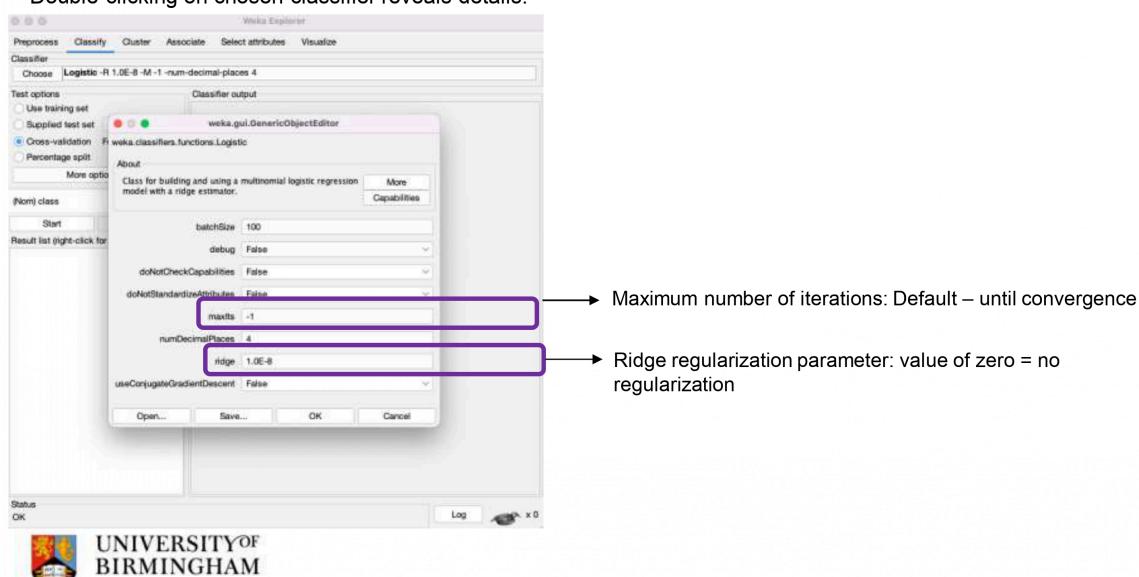
Example: Classification of Iris Data Set

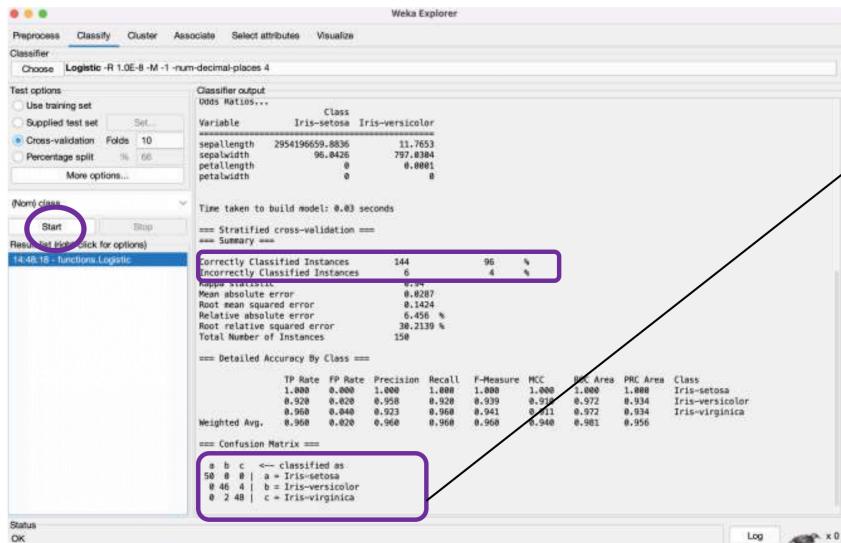


Logistic Regression



Double-clicking on chosen classifier reveals details.





Confusion matrix

- Each $(i,j)^{th}$ matrix entry corresponds to number of data instances where true class ' i ' is classified as class ' j '
- Row indicates true class labels, and column indicates predicted class labels

Confusion Matrix - Details

Confusion matrix is used to calculate class-specific quality metrics – accuracy, precision, recall, F-measure etc.

Simple binary classification (into apples and grapes)

- Corresponding to class apple (positive),

Predicted Values			
		POSITIVE (APPLE)	NEGATIVE (GRAPES)
Actual Values	POSITIVE (APPLE)	TP (True positive)	FN (False negative)
	NEGATIVE (GRAPES)	FP (False positive)	TN (True negative)

You predicted positive and it is true
You predicted negative and it is false
You predicted positive and it is false
You predicted negative and it is true

For more details: <https://towardsdatascience.com/understanding-confusion-matrix-a9ad42dcfd62>



Iris classification

- Corresponding to class a,

Predicted		
a	b	c
Actual	a	50 (cell 1) 0 (cell 2) 0 (cell 3)
b	0 (cell 4)	46 (cell 5) 4 (cell 6)
c	0 (cell 7)	2 (cell 8) 48 (cell 9)

TP = cell 1
(you predicted 'a' and it is true)

FN = cell 2+cell 3
(you predicted 'not a' and it is false)

FP = cell 4 +cell 7
(you predicted 'a' and it is false)

TN = cell 5+cell 6+ cell 8 + cell 9
(you predicted 'not a' and it is true)

Confusion Matrix - Details

Iris classification

- Corresponding to class b,

		Predicted		
		a	b	c
Actual	a	50 (cell 1)	0 (cell 2)	0 (cell 3)
	b	0 (cell 4)	46 (cell 5)	4 (cell 6)
c	0 (cell 7)	2 (cell 8)	48 (cell 9)	

TP = cell 5
FN = cell 4+cell 6
FP = cell 2 +cell 8
TN = cell 1+cell 3+ cell 7
+ cell 9

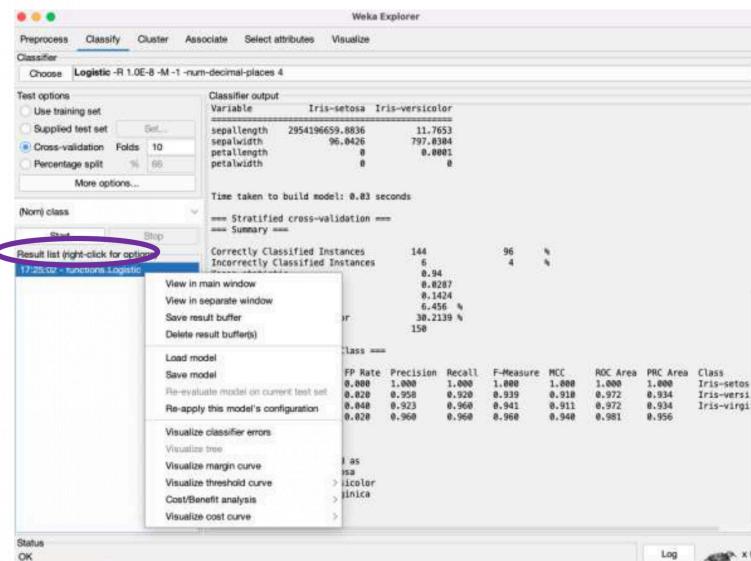
Per-class performance metrics:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad \text{Precision} = \frac{TP}{TP + FP} \quad \text{Recall} = \frac{TP}{TP + FN}$$

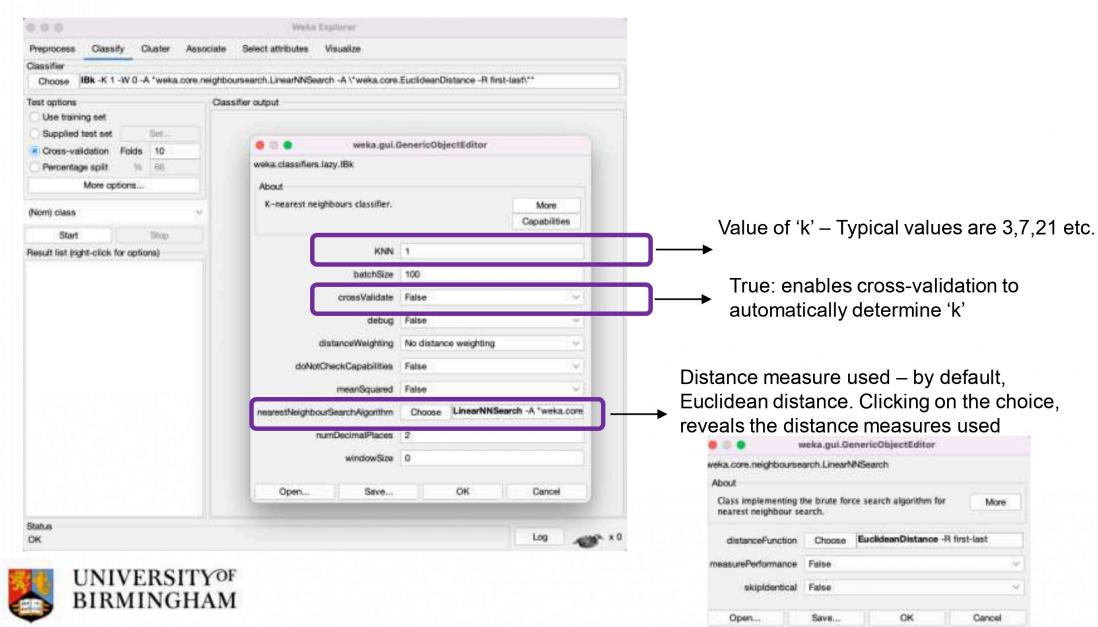
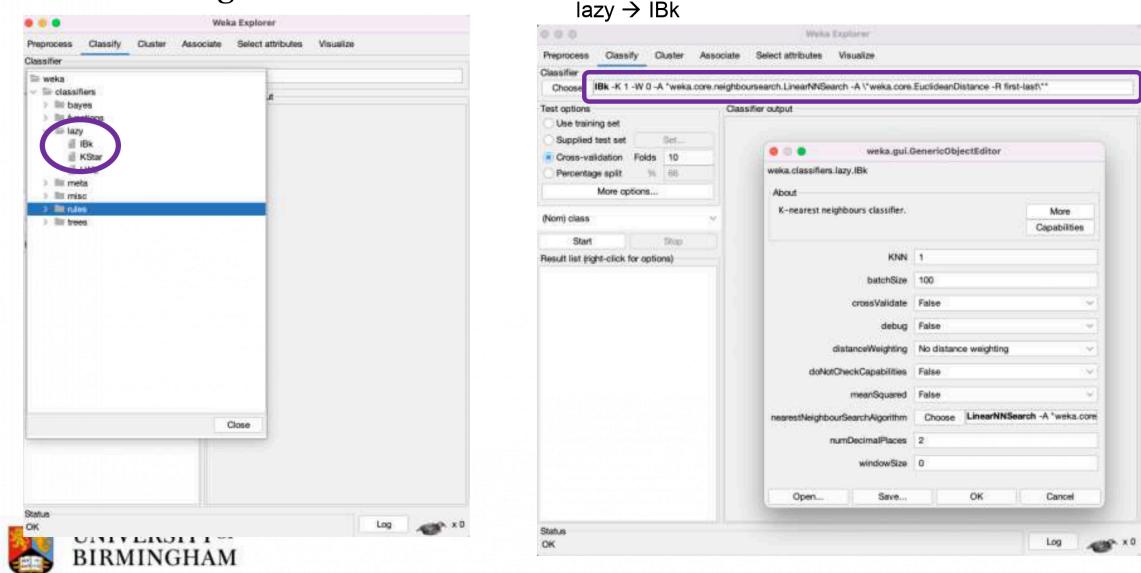
$$F_{\beta} = \frac{(\beta^2 + 1)Precision \times Recall}{\beta^2 Precision + Recall}$$

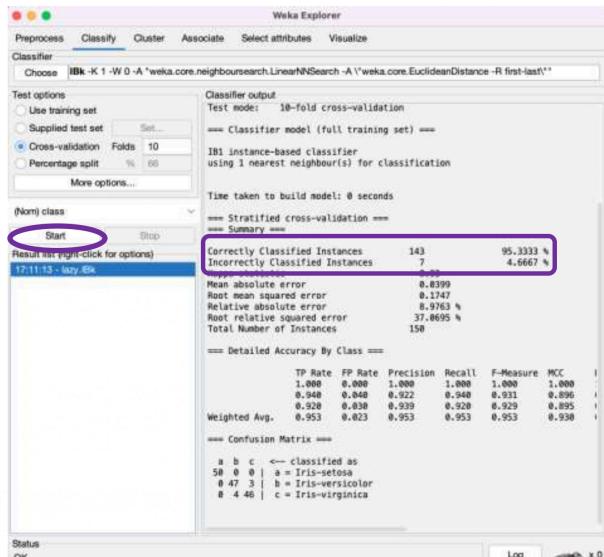
$\beta > 1$: Favor Recall
 $\beta < 1$: Favor Precision

$$\text{F-measure} = F_1$$

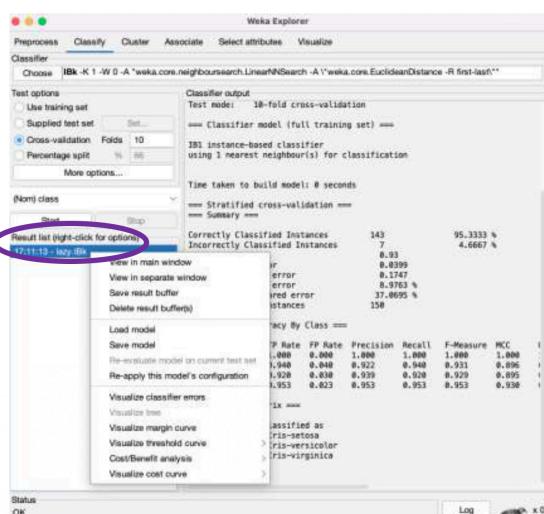


K-Nearest Neighbors





UNIVERSITY OF
BIRMINGHAM

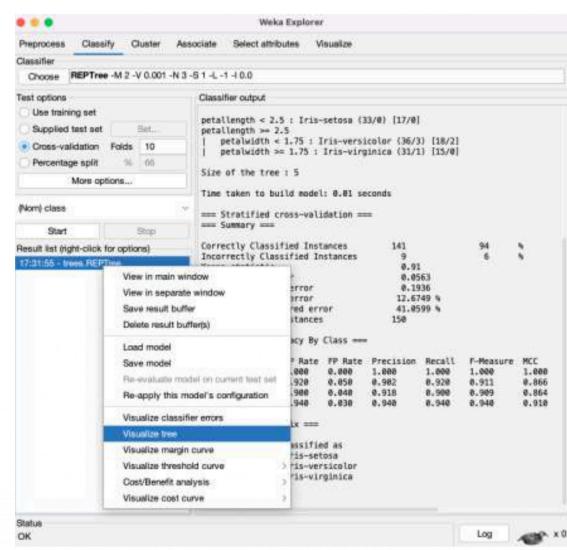
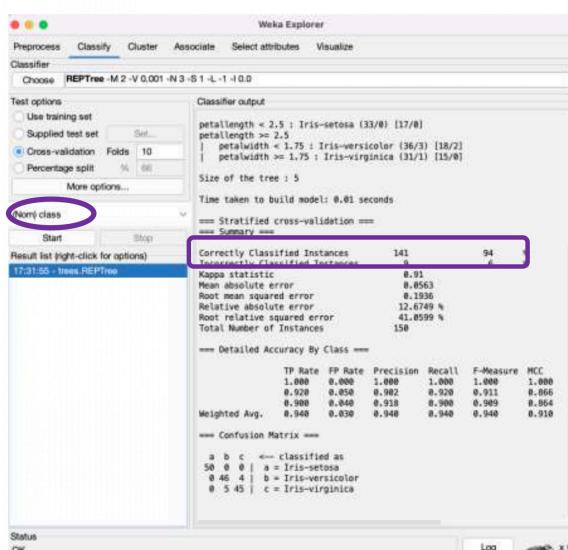
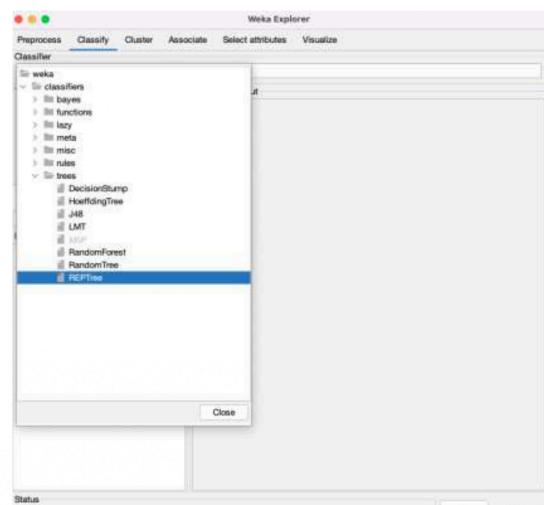


UNIVERSITY OF
BIRMINGHAM

Decision Trees

Decision Tree Algorithm:

- Tree is constructed in a top-down recursive divide-and-conquer manner
- At start, all the training examples are at the root
- 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)





ML Application: Sentiment Analysis



Sentiment Analysis: A classification task

- Sentiments:
 - Attitudes: positive/negative/neutral?
 - Emotions: happy/sad?
 - Opinions: like/dislike?
- Analysis of feeling behind the words
- Questions that might be asked in sentiment analysis:
 - Is the product review positive or negative?
 - Peoples' response to campaign reviews?
 - Bloggers' attitude about presidential election?



Sentiment Analysis with WEKA: A Route Map

- Get text data with labels (loading the .arff file).
- Pre-processing data: Convert text strings to feature vectors via NLP Techniques (use filters).
- Train a model using a ML algorithm for classification (use Classify → Logistic regression/ Decision Trees/k-NN).
- Test the model and use it for future classification



Example : Twitter dataset

- 100 positive tweets and 100 negative tweets from Edinburgh twitter corpus

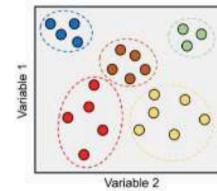
```
relation tweets
@attribute tweet_body string
@attribute sentiment {pos,neg}

@data
"i always feel motivated the fri afternoon prior to a holiday? wanted to get lots done... but i want
jammies and judge judy... \\"SIR!\\" @litj3 her ",pos
'seriously, do you have to rub it in maggie!!! ",pos
'if i could just get a job i would be employeed.. Tooltip is the \'title\' ",pos
'I don't like social karma much. Would rather skip it, but can't afford to piss my friend off
any more ",pos
'it's time to review the Ikeega if EMC send me one!! ",pos
'Something I have wanted to make for a while now... finally done URL",pos
'you're so sweet! proving me right again... the Dutch are the Best! ",pos
'last yr on the back of your Pepsi can! Our Pepsi can offer hit@torees today! Can't wait to
see one ",pos
'well rob i have to admit that you have to admit that you feel cool for being on twitter and
likeable to others ",pos
'big sky@ call at 1020! msg me if you want in ",pos
'well don't let it happen again ",pos
'reading the book about the world ",pos
'season 3... all over it... all its looking for season 4... i can't believe he didn't prepare... ",pos
'Yikes, it's windy here today; up to 58 mph NOAA says. Good thing electrician taught me generator
grill yesterday; Feeling uh, empowered. ",pos
'uh... i think it's a blue moon this weekend. It'll be a blue moon to-go. ",pos
'Makes at least two of us. ",pos
'that 'slang' doesn't really work for you ",pos
'yes i did THANK YOU ",pos
```

- Input: tweets, output: positive or negative
- Try in WEKA- Demo in lecture



UNIVERSITY OF
BIRMINGHAM



Explorer → Cluster

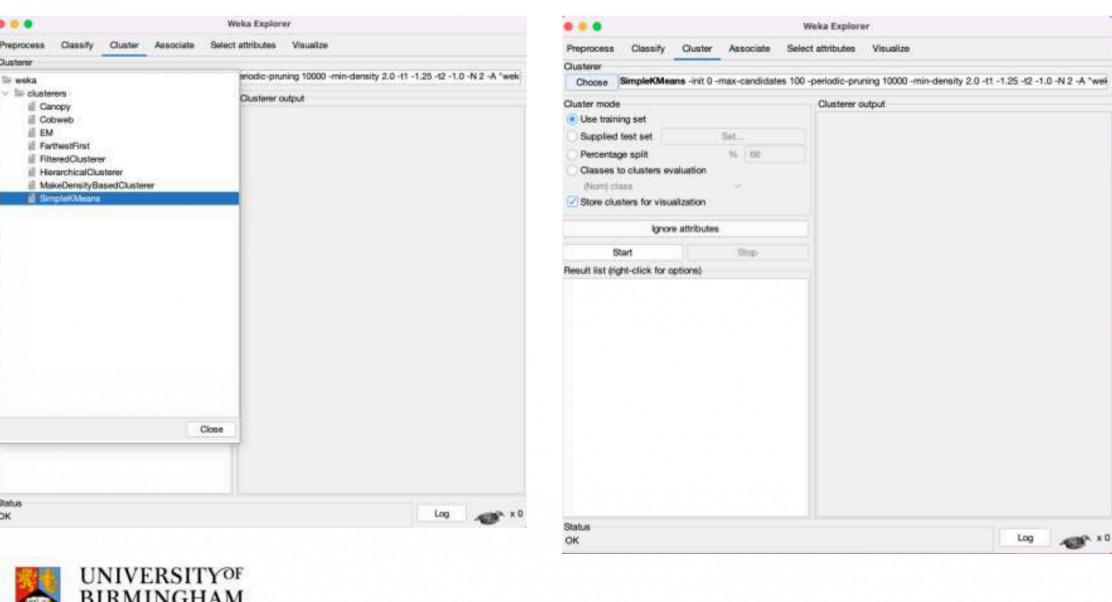
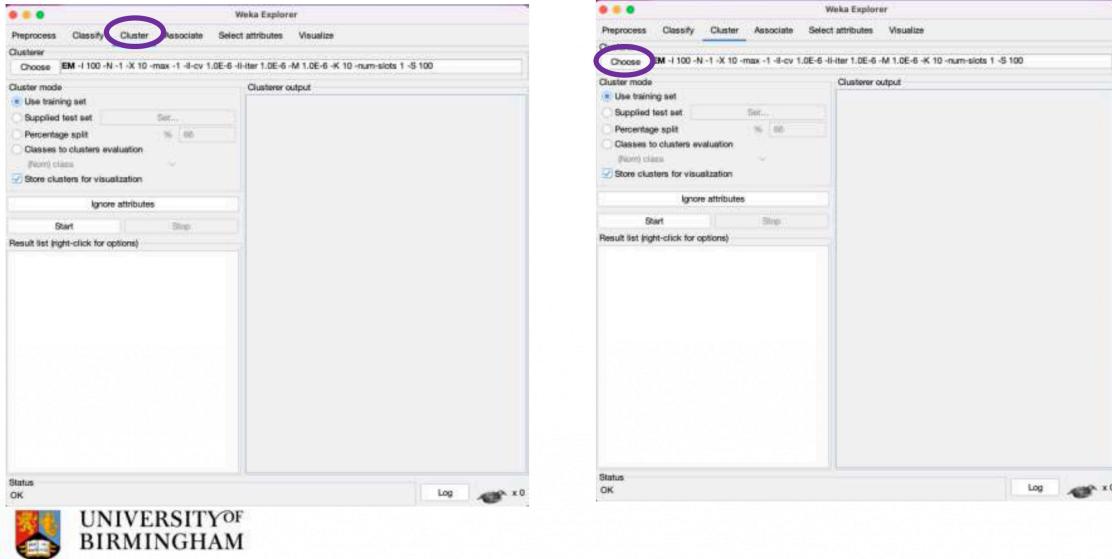
Clustering algorithms are a class of unsupervised learning algorithms (no labels) that aims to cluster the observed instances into groups with similar traits (to be covered in detail in the following lectures). WEKA supports clustering algorithms like:

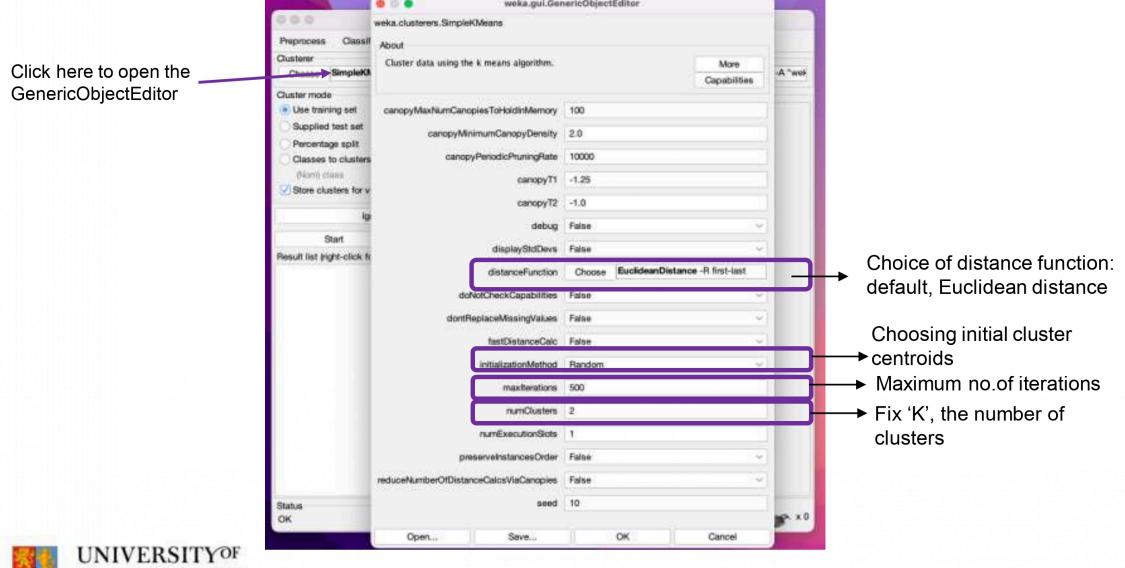
1. K-Means
2. Hierarchical Clustering
3. Density based clustering



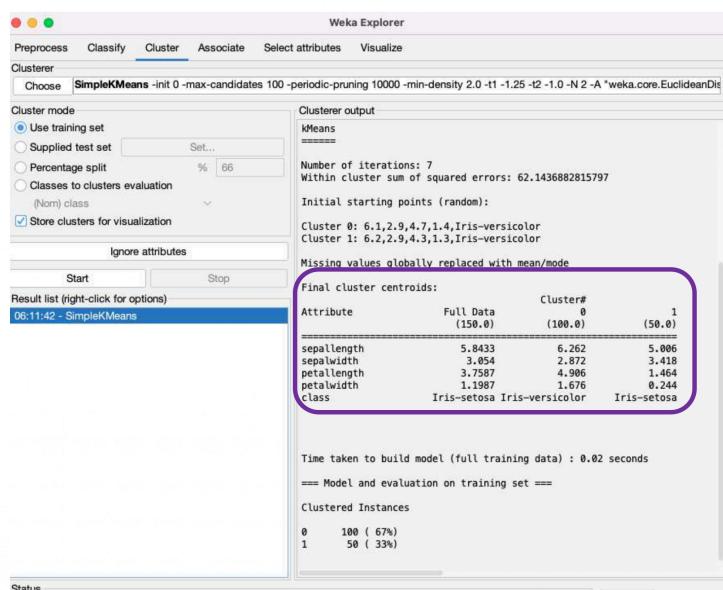
UNIVERSITY OF
BIRMINGHAM

Simple K-Means

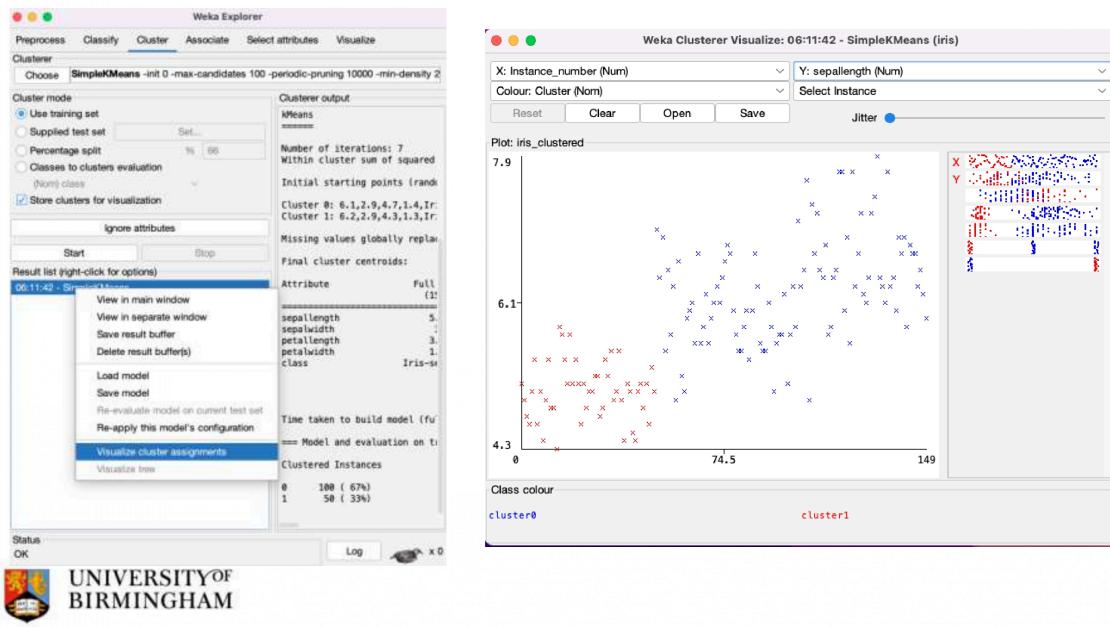




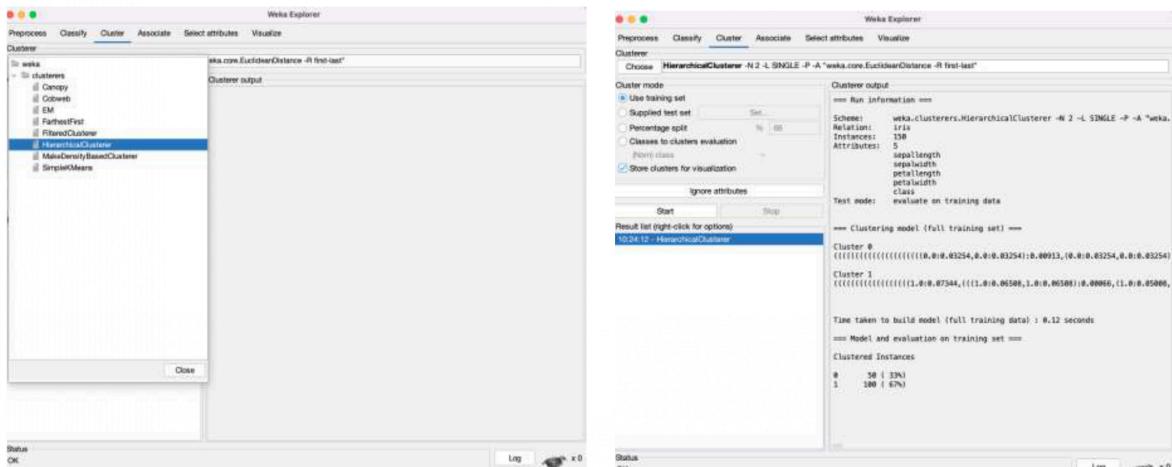
UNIVERSITY OF
BIRMINGHAM

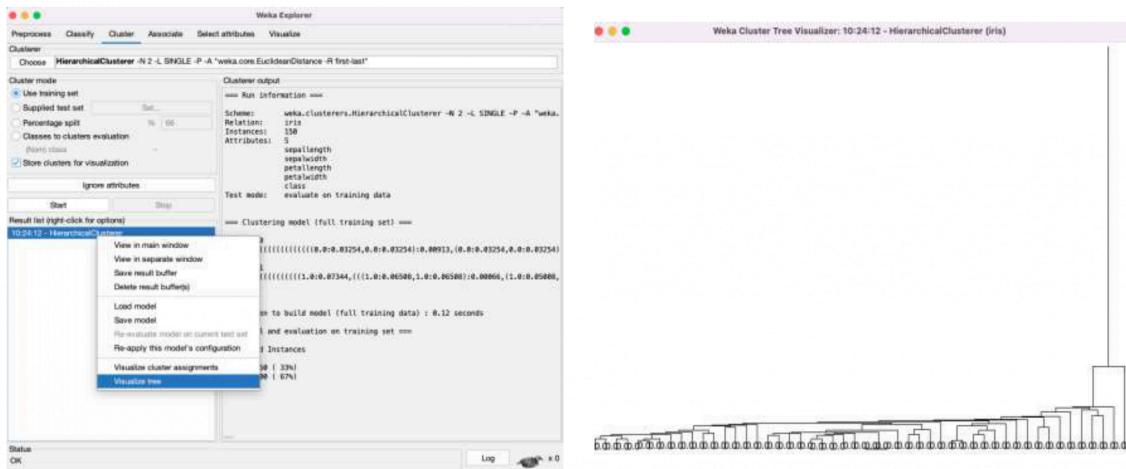


UNIVERSITY OF
BIRMINGHAM



Hierarchical Clusterer





UNIVERSITY OF
BIRMINGHAM

ML Application: Dominant Color Extraction



UNIVERSITY OF
BIRMINGHAM

Dominant Color Extraction – A Clustering Task

- Application – extraction of color palettes from artworks
- Input data – an image



```
@relation house
@attribute red real
@attribute green real
@attribute blue real

@data
26 20 45
28 5 46
28 12 44
28 13 46
31 4 51
31 5 48
31 8 44
31 9 44
32 5 50
32 8 47
32 11 47
32 12 48
33 9 47
33 17 45
33 29 48
34 6 48
34 9 50
34 13 52
34 28 45
34 32 47
```

Each data instance:
RGB value of each
pixel in the image



UNIVERSITY OF
BIRMINGHAM

Clustering Route Map in WEKA

- Load the .arff file corresponding to the RGB vectors of the image
- Train a model that cluster the RGB feature vectors into K groups via a clustering algorithm like K-means
- Obtain clusters and output the RGB value of each cluster centroid.

Note: cluster centroid = average of the feature vectors in a cluster (can be treated as a summary of the cluster)

Try in WEKA



UNIVERSITY OF
BIRMINGHAM

References

- For documentation on WEKA, <https://waikato.github.io/weka-wiki/documentation/>
- WEKA book: Ian H. Witten and Eibe Frank, Data Mining : Practical Machine Learning Tools and Techniques (Second Edition)

