



Introduction to WEKA
Preprocess and Classify

PL03



Material

<http://hpeixoto.github.io/dc>



WEKA



Weka is a collection of machine learning algorithms for data mining tasks.



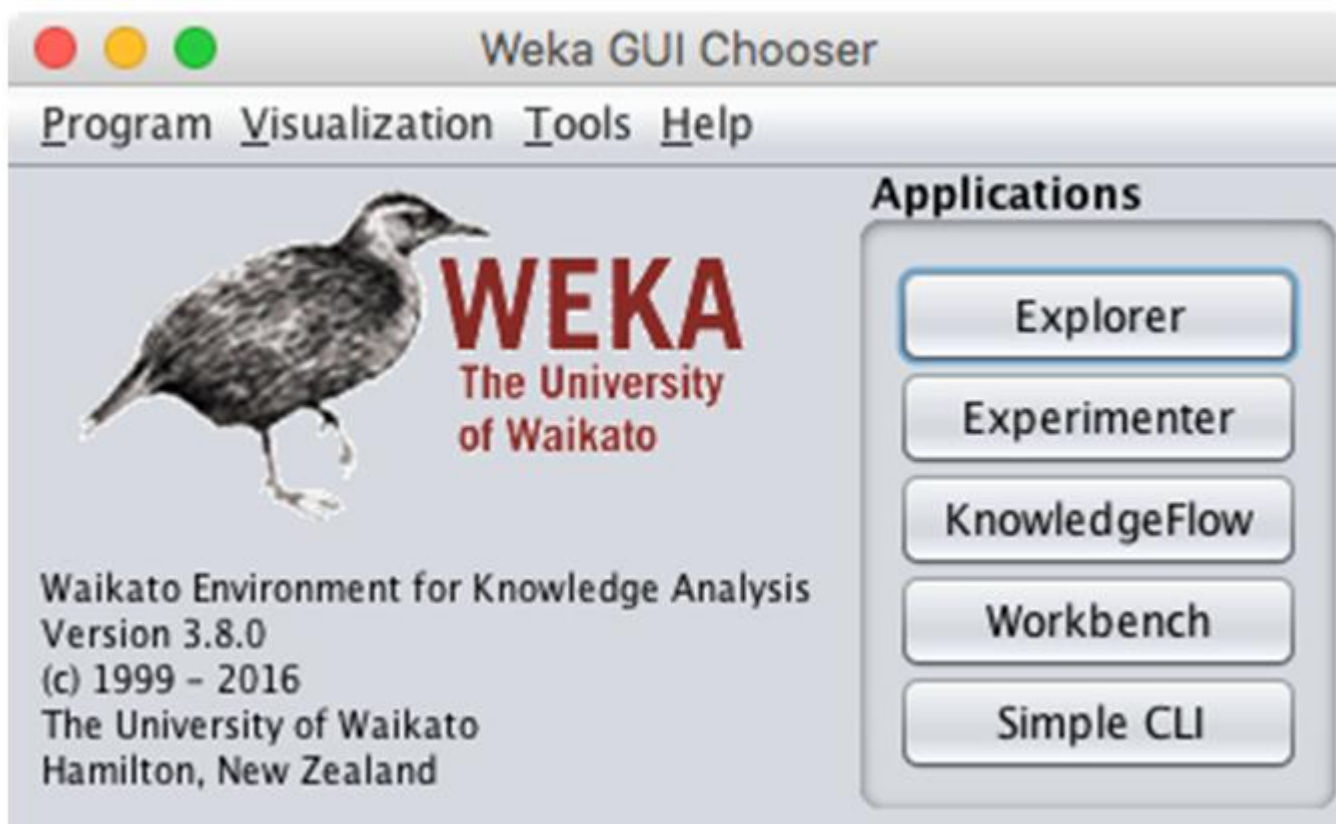
WEKA

WEKA download

<https://www.cs.waikato.ac.nz/ml/weka/downloading.html>

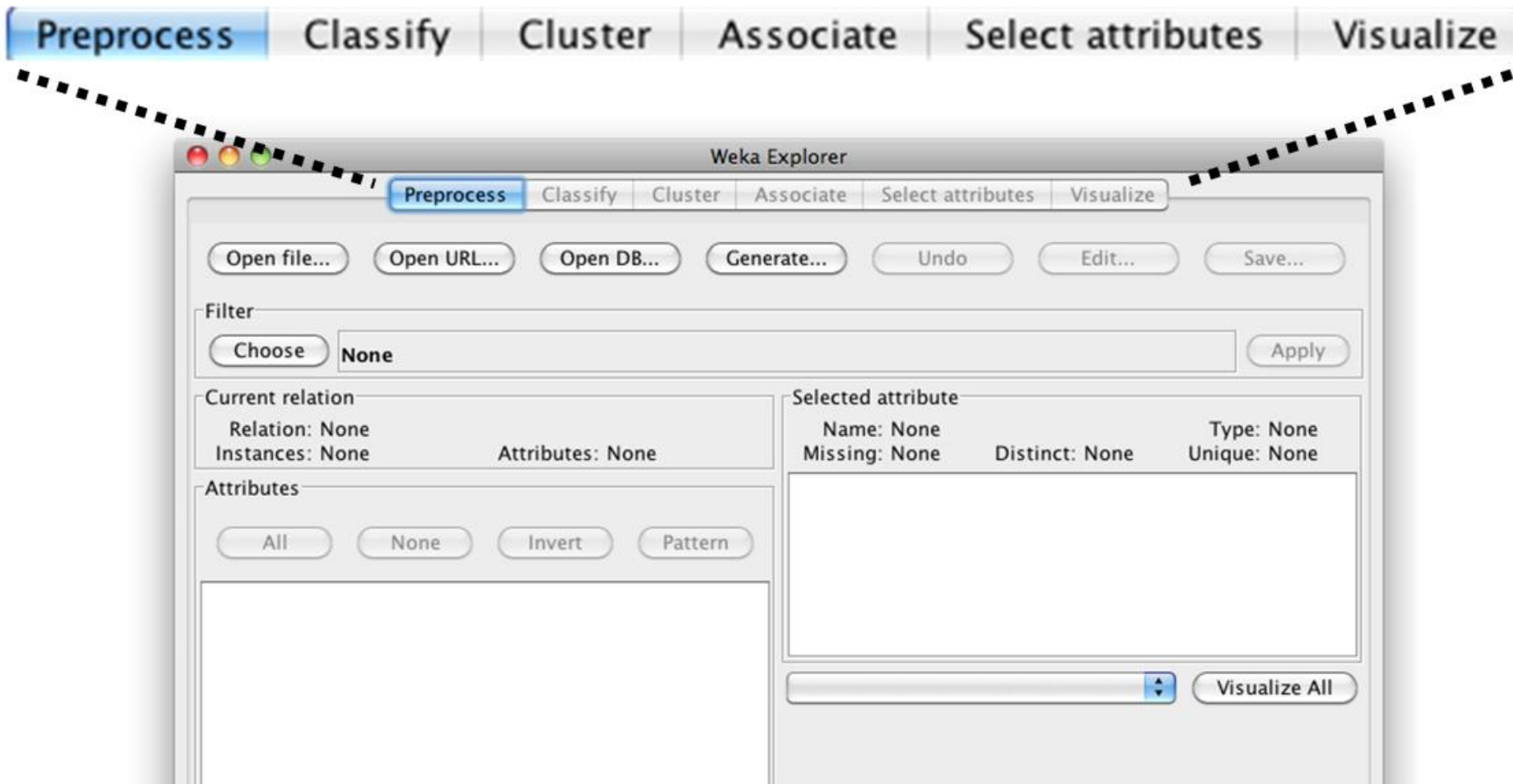


WEKA



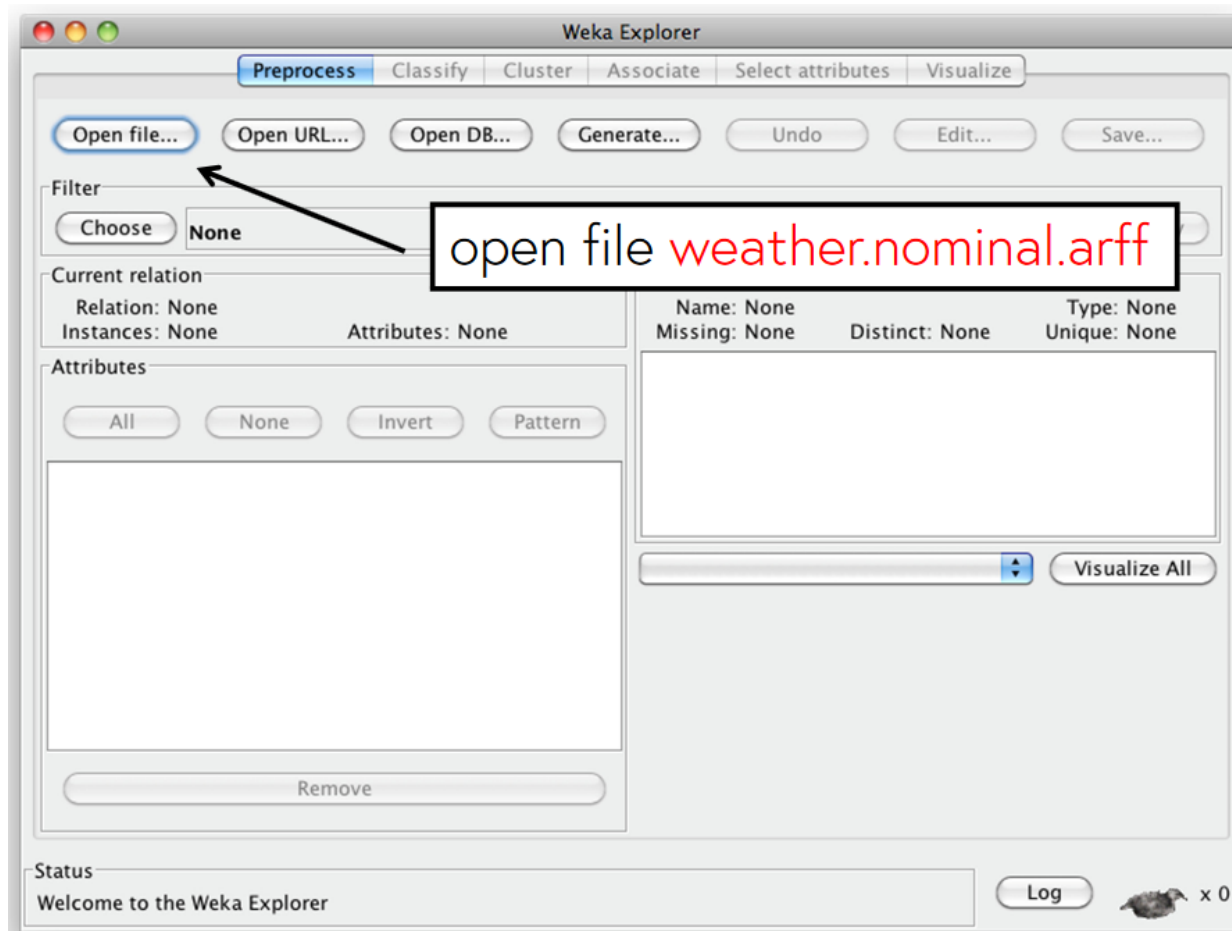


WEKA





WEKA





WEKA

The screenshot shows the Weka Explorer window with the 'Preprocess' tab selected. The 'Attributes' list on the left contains the following attributes:

No.	Name
1	outlook
2	temperature
3	humidity
4	windy
5	play

The 'Selected attribute' table on the right shows the distribution of the 'outlook' attribute:

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

The 'Class' is set to 'play (Nom)'. The 'Visualize All' button is visible. The status bar at the bottom shows 'OK' and a 'Log' button.



WEKA

Weather.arff

		attributes			
		Outlook	Temp	Humidity	Windy
instances	1	Sunny	Hot	High	False
	2	Sunny	Hot	High	True
	3	Overcast	Hot	High	False
	4	Rainy	Hot	High	False
	5	Rainy	Hot	Normal	False
	6	Rainy	Hot	Normal	True
	7	Rainy	Cool	Normal	True
	8	Rainy	Mild	Normal	False
	9	Sunny	Cool	Normal	False
	10	Rainy	Mild	Normal	False
	11	Sunny	Mild	Normal	True
	12	Overcast	Mild	High	True
	13	Overcast	Hot	Normal	False
	14	Rainy	Mild	High	True

Classification problem:
predict the "class" value

Play



WEKA

Weather.arff

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter: Choose None

Current relation: weather.symbolic
Instances: 14 Attributes: 5

Attributes: All None Invert Pattern

No.	Name
1	outlook
2	temperature
3	humidity
4	windy
5	play

Remove

Name: outlook
Missing: 0 (0%) Distinct: 3 Type: Nominal
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

class

attribute values

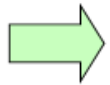


WEKA

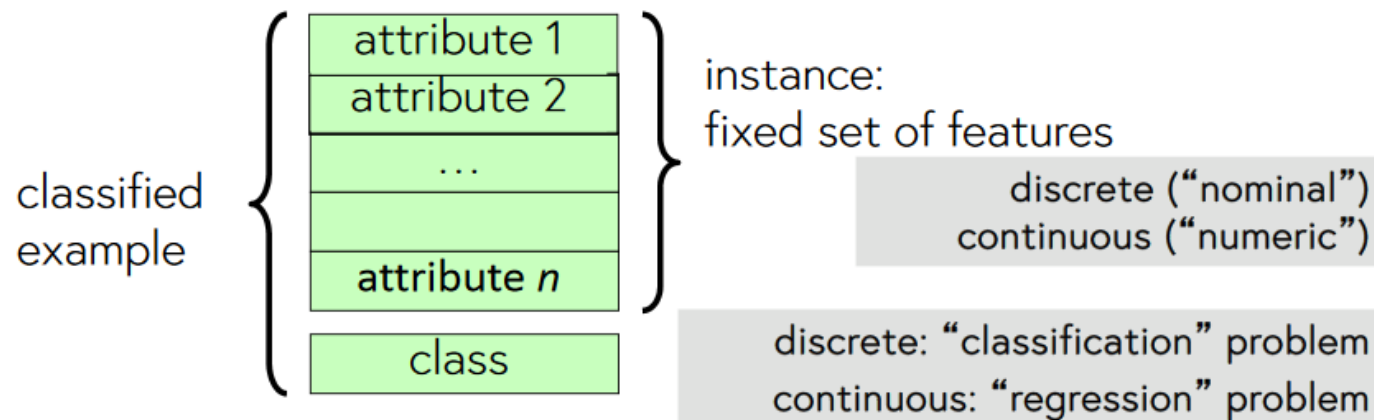
Classification

sometimes called “supervised learning”

Dataset: classified examples



“Model” that classifies new examples





WEKA

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter: Choose None

Current relation: Relation: weather.symbolic Instances: 14 Attributes: 5

Attributes: All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> outlook
2	<input type="checkbox"/> temperature
3	<input type="checkbox"/> humidity
4	<input type="checkbox"/> windy
5	<input type="checkbox"/> play

Remove

Status: OK

Log x 0

open file weather.numeric.arff

attribute values

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

Class: play (Nom) Visualize All



BUILD CLASSIFIER

USE J48 TO ANALYZE THE GLASS DATASET

- Open file `glass.arff`
- 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



BUILD 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



BUILD CLASSIFIER

Pruning (decision trees)

is a technique in machine learning that reduces the size of **decision trees** by removing sections of the **tree** that provide little power to classify instances. **Pruning** reduces the complexity of the final classifier, and hence improves predictive accuracy by the reduction of overfitting.



EXERCÍCIO GRUPO

DETERMINAR NECESSIDADE DE COMPONENTES SANGUÍNEOS

No bloco operatório de uma unidade de saúde é de vital importância determinar com a devida antecedência a potencial necessidade de um paciente, que vai ser intervencionado, vir a receber uma ou mais transfusões de componentes sanguíneos (sangue, plasma, etc).

Esta necessidade advém de diversos fatores, tais como:

- Custo unitário de cada fornecimento de componentes sanguíneos;
- Escassez de oferta de componentes sanguíneos;
- Correto tratamento do paciente em caso de necessidade de transfusão;
- Diminuir os desperdícios com o deteriorar de componentes sanguíneos afetos a cirurgias que não são usados;
- entre outros...



EXERCÍCIO

DETERMINAR NECESSIDADE DE COMPONENTES SANGUÍNEOS

Num exercício de brainstorm refletir sobre os principais dados que deveriam estar acessíveis.

nota: Não existe nenhuma limitação técnica nem legal para a recolha de qualquer input. Pelo que todos os dados estarão acessíveis caso sejam considerados.



EXERCISES

FE02



Introduction to WEKA
Preprocess and Classify

PL03