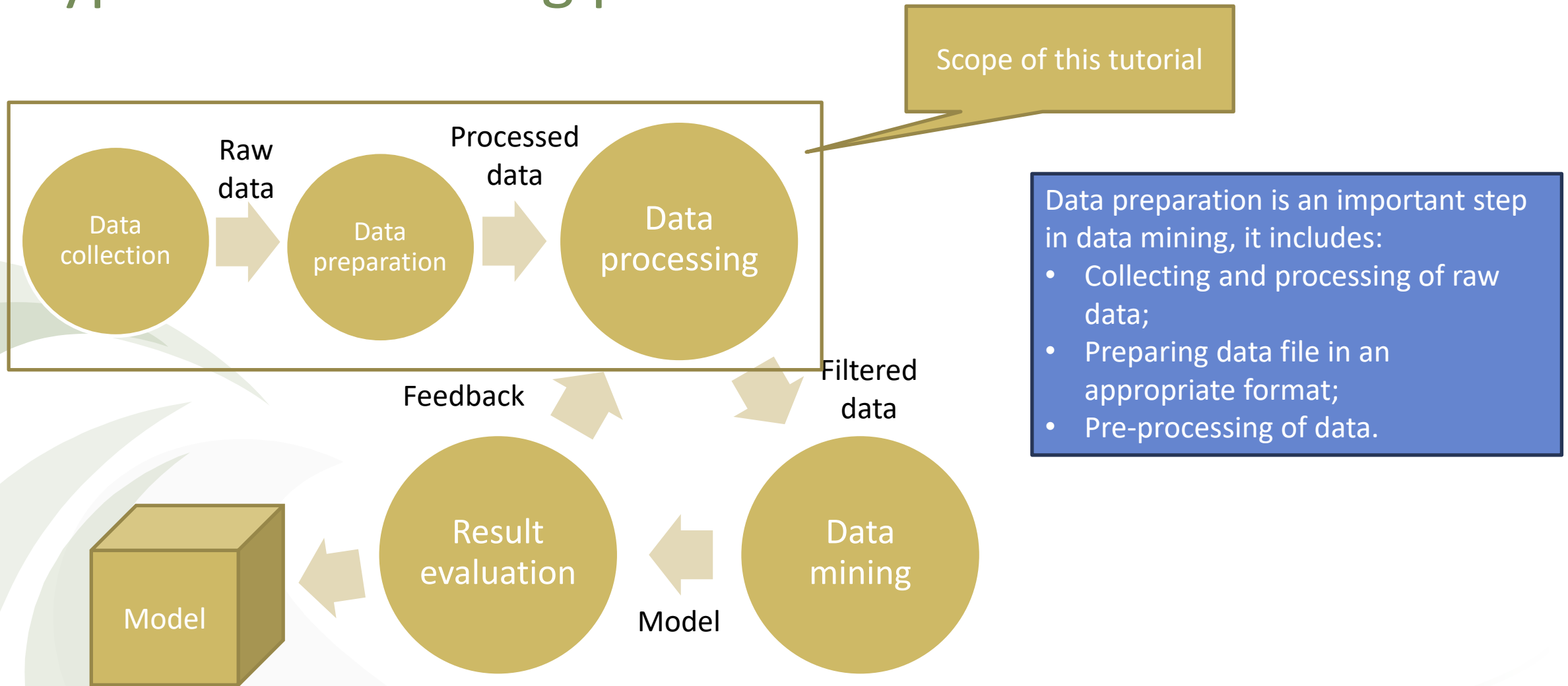# COMP7103 Data Mining

**Tutorial 1**

Data Preparation / Weka

# Typical data mining process

# Data preparation tools

An open-source data mining software in Java
https://www.cs.waikato.ac.nz/ml/weka/
(Download version 3.8 or above)

- **Weka**
  - Allow you to apply filters to existing data set
  - Not useful if raw data is in a format not readable by Weka

- **Spreadsheets (Excel, Google Sheets, etc.)**
  - Good for small data set, especially if you are familiar with the formulas
  - Explore potion of dataset before action
  - Cannot process large amount of data

- **Own program**
  - Require the most effort
  - Most customizable

In any cases, data needs to be converted into a format appropriate for the data mining exercise depending on the choice of tools.

# Example dataset for this tutorial

- 150 Iris data is collected in a CSV file
  - https://archive.ics.uci.edu/ml/datasets/iris

- There are 4 attributes
  - *Sepal* **length** and **width**
  - *Petal* **length** and **width**

- Three species of iris

> Available on Moodle:
> `iris.csv`


Iris Setosa


Iris Versicolor


Iris Virginica

# Understanding the data

- Read carefully the information in the data source.
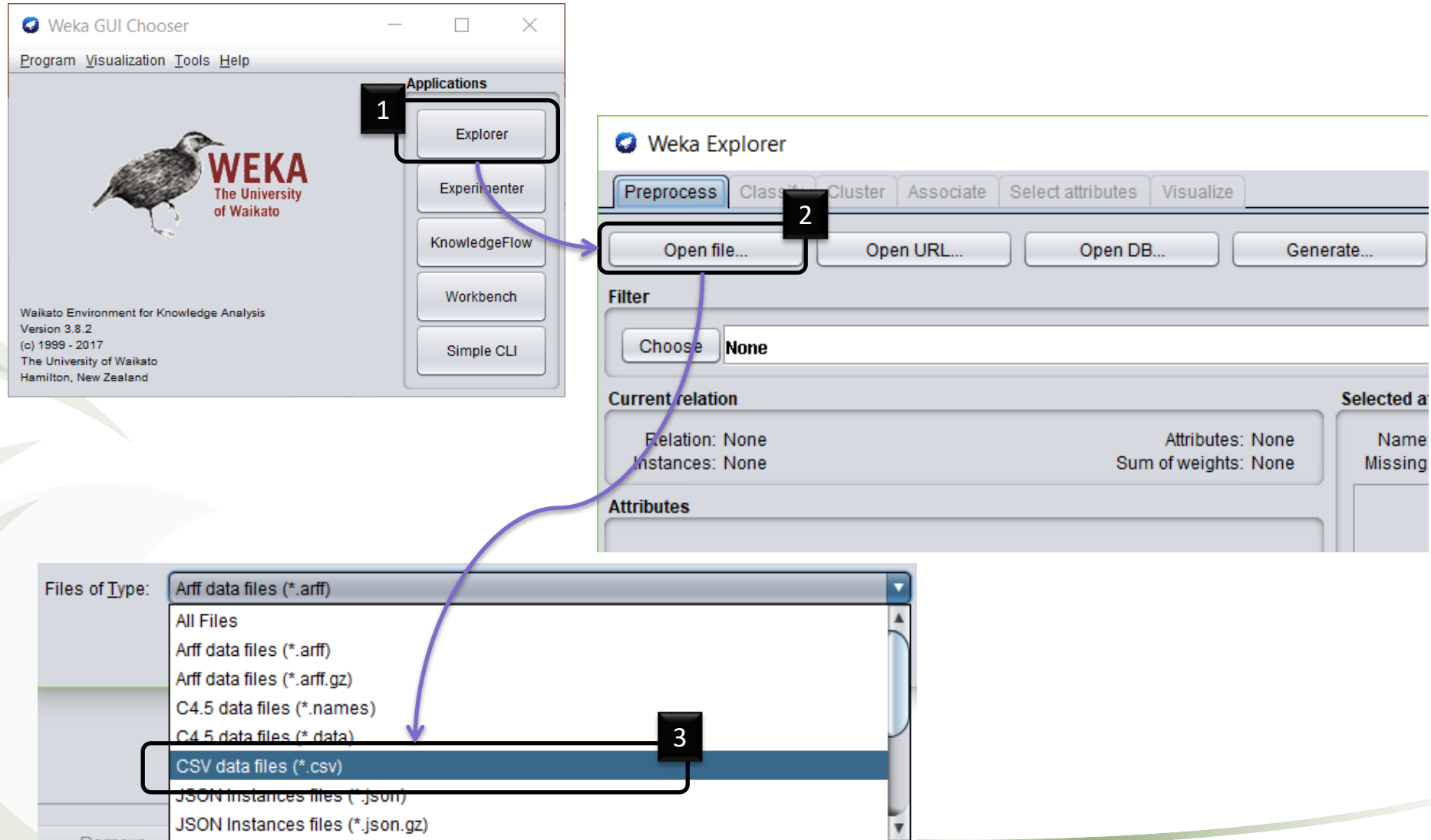    - https://archive.ics.uci.edu/ml/datasets/iris

> "…The 35th sample should be: 4.9,3.1,1.5,0.2,"Iris-setosa" where the error is in the fourth feature. The 38th sample: 4.9,3.6,1.4,0.1,"Iris-setosa" where the errors are in the second and third features."

- Plan the preprocessing steps:
    - Fix data as described in the data source
    - Remove attributes that may not be useful
    - Save data for future use

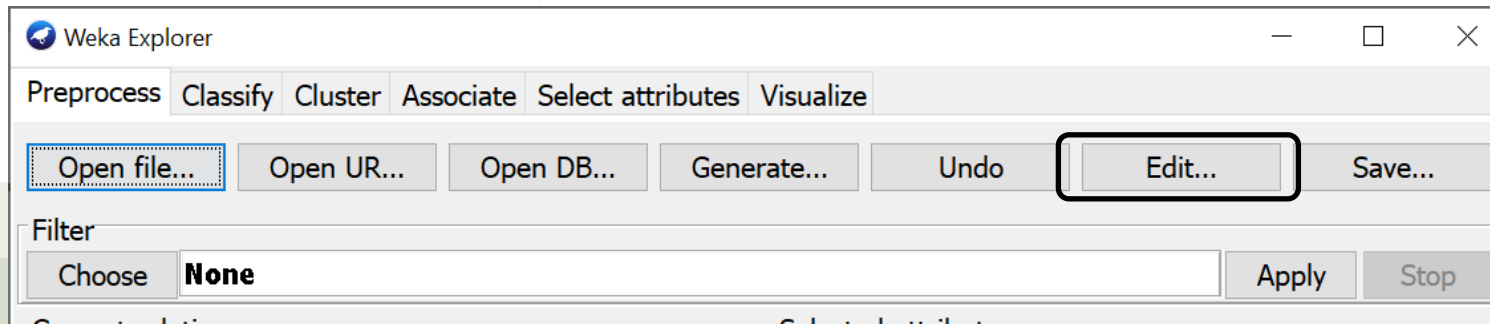Pre-processing CSV file in Weka

# Data preparation

# Weka Explorer – opening CSV file

# Edit data

- Open the data in Weka and edit the data.



Edit these values accordingly

# Remove attribute (if needed)



ID is not useful in classification, we can remove it

# Save data

- Then we can save the data into an ARFF file

# ARFF file

```
% comments
@RELATION relation_name

@ATTRIBUTE attribute_name attribute_type
@ATTRIBUTE attribute_name attribute_type
...

@DATA
comma-separated values
```

- An ARFF file is a plain-text file with a specific format:
  - https://waikato.github.io/weka-wiki/formats_and_processing/arff_stable/
- For example, this is the iris dataset:

Data-type can be either:
- numeric
- {list,of,normal,values}
- string
- date [<date-format>]

```
% Iris Plants Database
@RELATION iris

@ATTRIBUTE sepallength NUMERIC
@ATTRIBUTE sepalwidth  NUMERIC
@ATTRIBUTE petallength NUMERIC
@ATTRIBUTE petalwidth  NUMERIC
@ATTRIBUTE class       {Iris-setosa,Iris-versicolor,Iros-virginica}

@DATA
5.1,3.5,1.4,0.2,Iris-setosa
4.9,3.0,1.4,0.2,Iris-setosa
...
```

After this line, the data is presented in CSV format

11

# Weka Explorer – opening ARFF file

Then we can open the ARFF file using Weka

# Visualize one attribute

| Maximum | 7.9 |
|---------|-----|
| Mean | 5.843 |
| StdDev | 0.828 |

**1** Class: class (Nom)

**2** Visualize All

Petal Length looks very significant in predicting one of the class

# Using plot matrix

# Interpreting the matrix



Looks good: classes are separated

Looks poor: classes are mixed together

Is this the only way to interpret the matrix?

# Interpreting the matrix again



Not good, attributes are correlated

Seems good: attributes not correlated

There could be many ways to interpret the data!

16

# Applying filter

- Applying filter under the "Pre-process" tab of Weka is one of the easiest way to preprocess the data.

# Result

Data Denormalization

# Pre-processing Raw data

# Pre-processing raw data

- In this demonstration, we are using data from the UCI Machine Learning Repository.
  - https://archive.ics.uci.edu/ml/datasets/Anonymous+Microsoft+Web+Data

```
I,4,"www.microsoft.com","created by getlog.pl"
T,1,"VRoot",0,0,"VRoot"
N,0,"0"
N,1,"1"
T,2,"Hide1",0,0,"Hide"
N,0,"0"
N,1,"1"
A,1287,1,"International AutoRoute","/autoroute"
A,1288,1,"library","/library"
A,1289,1,"Master Chef Product Information","/masterchef"
A,1297,1,"Central America","/centroam"
A,1215,1,"For Developers Only Info","/developer"
...
```

Available on Moodle: `webdata.csv`

# Understanding the raw data

- The raw data is in CSV form

- For each row, the first column indicate the type of the data:
  - 'A' for attribute (a page);
  - 'C' for case (a user);
  - 'V' for vote (a visit of a page);
  - All others are ignored in our case.

- So basically, we can create a set of relational data from the data set:

```
...
A,1008,1,"Free Downloads","/msdownload"
...
A,1046,1,"IE Support","/iesupport"
...
A,1034,1,"Internet Explorer","/ie"
...
C,10027,10027
V,1008,1
V,1046,1
V,1034,1
...
```

| user_id |
|---------|
| … |
| 10027 |
| … |

**Users**

| user_id | page_id |
|---------|---------|
| 10027 | 1008 |
| 10027 | 1034 |
| 10027 | 1046 |

**Visits**

| page_id | path |
|---------|------|
| 1008 | /msdownload |
| 1034 | /ie |
| 1046 | /iesupport |

**Pages**

# What do we need?

| Users | Visits | Pages |
|---|---|---|
| *user_id* | *user_id, page_id* | *page_id, url* |

- What is the purpose of our data mining exercise?

- What is the data mining exercise that we are going to do?

- Suppose, we only focus on classifying whether a user will visit a page, we need a data set like this:

| User_id | page 1 | page 2 | page 3 | page 4 | page 5 | ... |
|---|---|---|---|---|---|---|
| 1 | Yes | No | Yes | No | No | ... |
| 2 | No | No | Yes | No | No | ... |
| ... | | | | | | |

What can we do?

# Denormalization

| user_id | ... | 1008 | 1034 | 1046 | ... |
|---------|-----|------|------|------|-----|
| ... | ... | ... | ... | ... | ... |
| 10027 | ... | Yes | Yes | Yes | ... |
| ... | | | | | |

- We can denormalize the data to fit our purpose.

```
...
A,1008,1,"Free Downloads","/msdownload"
...
A,1046,1,"IE Support","/iesupport"
...
A,1034,1,"Internet Explorer","/ie"
...
C,10027,10027
V,1008,1
V,1046,1
V,1034,1
...
```

| user_id | page_id | path |
|---------|---------|------|
| 10027 | 1008 | /msdownload |
| 10027 | 1034 | /ie |
| 10027 | 1046 | /iesupport |

| user_id |
|---------|
| ... |
| 10027 |
| ... |

**Users**

| user_id | page_id |
|---------|---------|
| 10027 | 1008 |
| 10027 | 1034 |
| 10027 | 1046 |

**Visits**

| page_id | path |
|---------|------|
| 1008 | /msdownload |
| 1034 | /ie |
| 1046 | /iesupport |

**Pages**

# Creating `.arff` file for Weka

- By implementing a custom program to process the raw data, a CSV file is generated:

```
activeplatform,activex,athome,corpinfo,education,exchange, …
N,N,Y,N,N,N,...
N,N,N,N,N,N,...
N,N,N,N,N,N,...
```

Available on Moodle:
**webdata.processed.csv**

- One can easily convert it to `.arff` file by converting the first row (the labels) into ARFF header.

```
@relation web_log_data

@attribute activeplatform {Y,N}
@attribute activex         {Y,N}
...

@data
N,N,Y,N,N,N,...
N,N,N,N,N,N,...
...
```

# Auto-conversion by Weka

- If you open the CSV file in Weka and save it as ARFF file, you may need to check the nominal attributes for compatibility.

```
@relation web.data.processed

@attribute activeplatform {N,Y}
@attribute activex {N,Y}
@attribute athome {Y,N}
...

@data
N,N,Y,N,N,N,...
N,N,N,N,N,N,...
...
```

The two attributes will be considered incompatible in Weka