

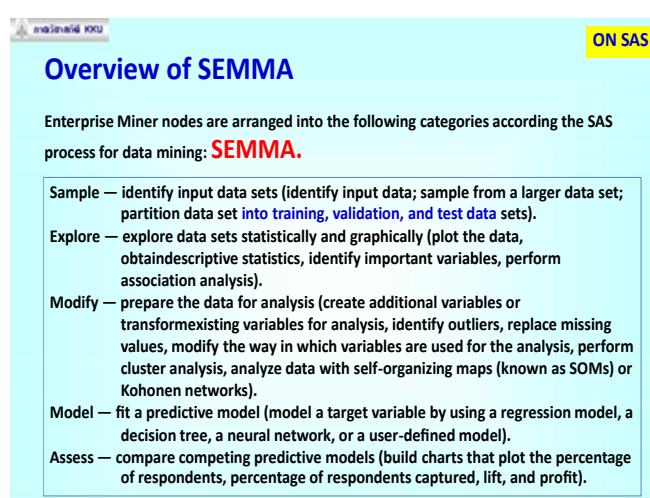
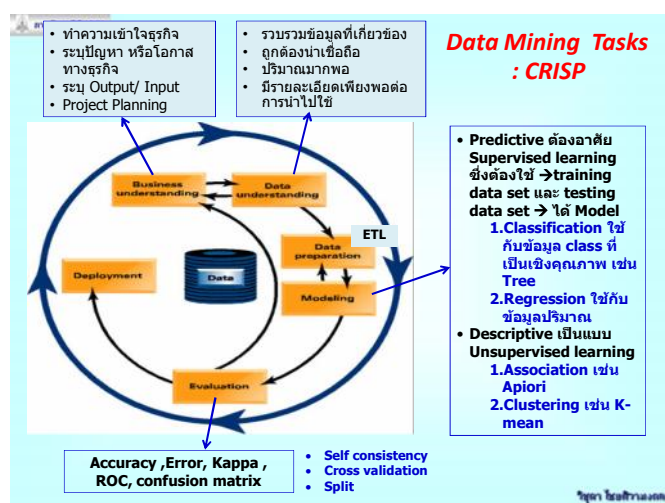
## การใช้โปรแกรม WEKA

**Weka** (Waikato Environment for Knowledge Analysis) is a collection of machine learning algorithms for data mining tasks. The algorithms can either be applied directly to a dataset or called from your own Java code. Weka contains tools for data pre-processing, classification, regression, clustering, association rules, and visualization. It is also well-suited for developing new machine learning schemes.

The Explorer interface features several panels providing access to the main components of the workbench:

- **The Preprocess panel** has facilities for importing data from a database, a CSV file, etc., and for preprocessing this data using a so-called filtering algorithm. These filters can be used to transform the data (e.g., turning numeric attributes into discrete ones) and make it possible to delete instances and attributes according to specific criteria.
- **The Classify panel** enables the user to apply classification and regression algorithms (indiscriminately called classifiers in Weka) to the resulting dataset, to estimate the accuracy of the resulting predictive model, and to visualize erroneous predictions, ROC curves, etc., or the model itself (if the model is amenable to visualization like, e.g., a decision tree).
- **The Associate panel** provides access to association rule learners that attempt to identify all important interrelationships between attributes in the data.
- **The Cluster panel** gives access to the clustering techniques in Weka, e.g., the simple k-means algorithm. There is also an implementation of the expectation maximization algorithm for learning a mixture of normal distributions.
- **The Select attributes panel** provides algorithms for identifying the most predictive attributes in a dataset.
- **The Visualize panel** shows a scatter plot matrix, where individual scatter plots can be selected and enlarged, and analyzed further using various selection operators.

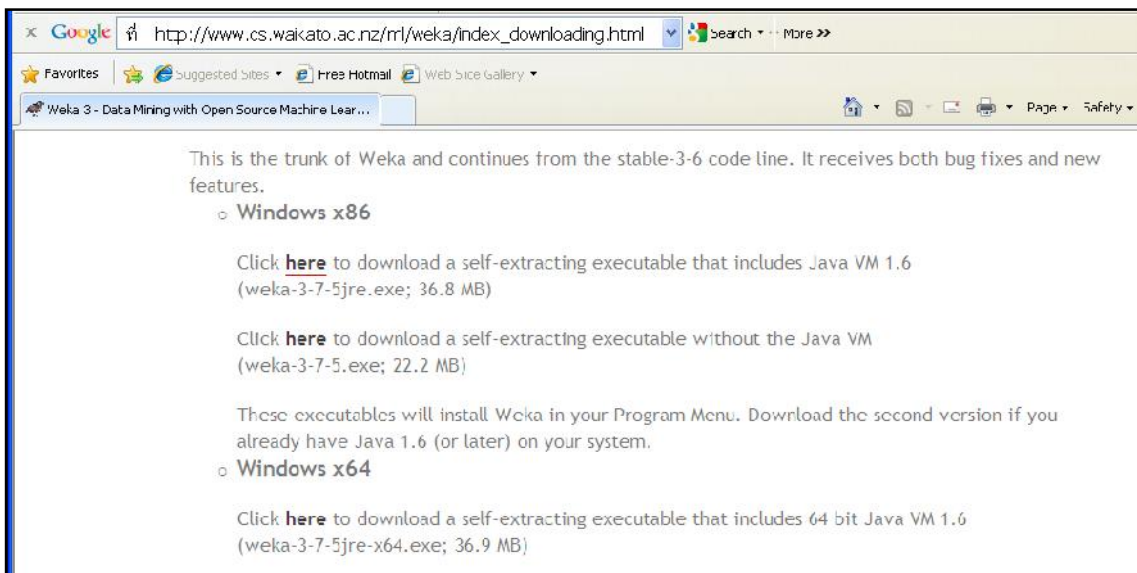
โปรแกรม Weka เริ่มพัฒนามาตั้งแต่ปี 1997 โดยมหาวิทยาลัย Waikato ประเทศนิวซีแลนด์ เป็นซอฟต์แวร์สำเร็จรูปประกอบประเภทฟรีแวร์ อยู่ภายใต้การควบคุมของ GPL License ซึ่งโปรแกรม Weka ได้ถูกพัฒนามาจากภาษาจาวาทั้งหมด ซึ่งเขียนมาโดยเน้นกับงานทางด้านการศึกษาด้วยเครื่อง (Machine Learning) และ การทำเหมืองข้อมูล (Data Mining) โปรแกรมจะประกอบไปด้วยโมดูลย่อย ๆ สำหรับใช้ในการจัดการข้อมูล และเป็นโปรแกรมที่สามารถใช้ Graphic User Interface (GUI) โดยมีฟังก์ชันสำหรับการทำงานร่วมกับข้อมูล ได้แก่ Pre-Processing, Classification, Regression, Clustering, Association rules, Selection และ Visualization



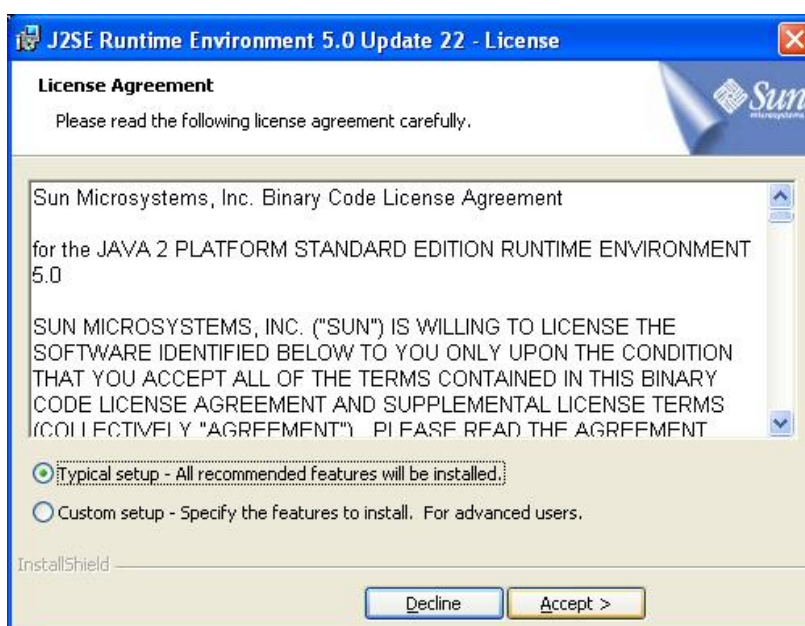
## การ Install Software

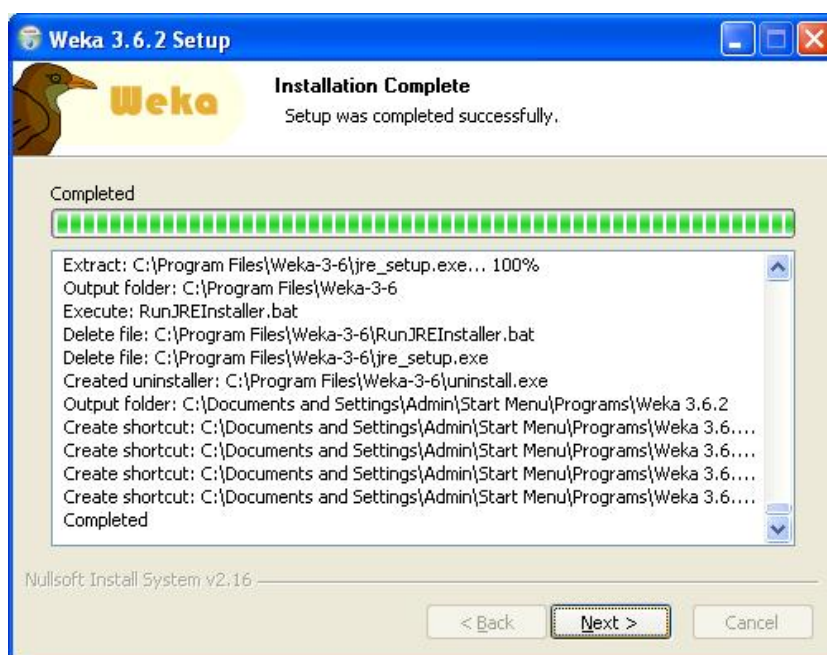
กรณีที่ในเครื่องท่านยังไม่มี **JAVA** ให้ download file แบบ **includes JAVA** โดยสามารถเข้าที่ website ต่อไปนี้

[http://www.cs.waikato.ac.nz/ml/weka/index\\_downloading.html](http://www.cs.waikato.ac.nz/ml/weka/index_downloading.html)

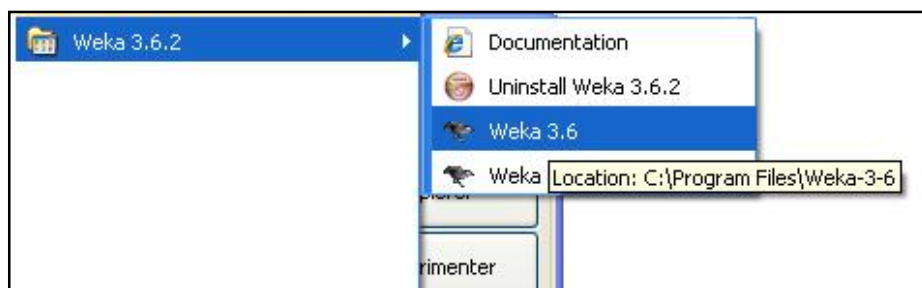


จากนั้นทำการ run



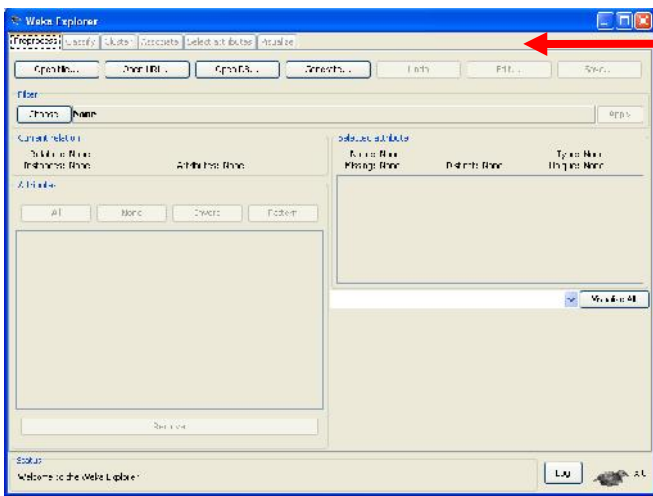


## การเรียกใช้



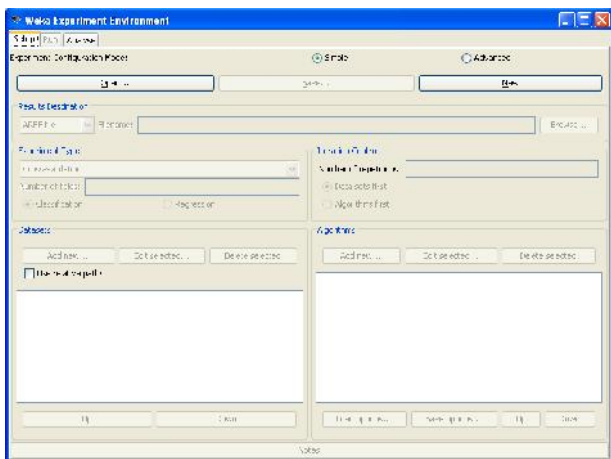


- Explorer เป็นโปรแกรมที่ออกแบบในลักษณะ GUI
- Experimenter เป็นโปรแกรมที่ออกแบบการทดลองและการทดสอบผล
- KnowledgeFlow เป็นโปรแกรมออกแบบผังการไหลของความรู้
- Simple CLI (Command Line Interface) เป็นโปรแกรมรับคำสั่งการทำงานผ่านการพิมพ์

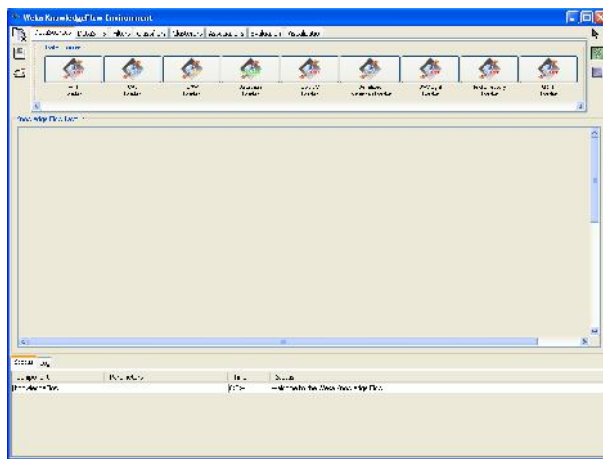


Explorer

- **Preprocess** การเตรียมข้อมูล เลือก file input พิจารณารายละเอียดข้อมูล แก้ไขข้อมูล แปลงข้อมูล
- **Classify** เป็นโมดูลการทำเหมืองข้อมูลแบบการจำแนกประเภท (Classification) จำแนกประเภทข้อมูลทำนายค่าข้อมูลใหม่จาก train model
- **Cluster** เป็นโมดูลการทำเหมืองข้อมูลแบบการแบ่งกลุ่ม (Clustering) แบ่งกลุ่มข้อมูลตามความคล้ายคลึง (Similarity)
- **Associate** เป็นโมดูลการทำเหมืองข้อมูลแบบกฎความสัมพันธ์ (Association rule)
- **Select attribute** คัดเลือกตัวแปรที่สำคัญ
- **Visualize** แสดงผลของข้อมูลในรูปแบบต่างๆ สองมิติ



Experimenter



KnowledgeFlow

## การนำข้อมูลเข้าสู่ WEKA

### ประเภทของแฟ้มข้อมูลที่ได้รับได้

- แฟ้มข้อมูลที่ได้รับต้องเป็น ARFF หรือ CSV
- ในกรณีที่แฟ้มข้อมูลอยู่ในเครือข่ายสามารถเรียกใช้ผ่าน URL ได้
- สามารถเรียกใช้ข้อมูลจากฐานข้อมูลได้ โดยเชื่อมโยงผ่าน JDBC

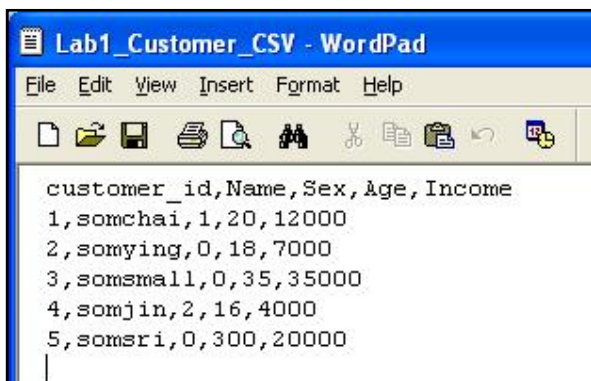
กรณี **Attribute Relationship File Format (ARFF)** is the text format file used by Weka to store data in a database

```
@relation weather
@attribute outlook {sunny, overcast, rainy}
@attribute temperature real
@attribute humidity real
@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
```

- **@relation name** เป็นบรรทัดที่บอกชื่อตารางข้อมูลเชิงสัมพันธ์
- **@attribute att-name type** เป็นบรรทัดที่บอกชื่อลักษณะเฉพาะและชนิด
  - numeric หรือ real หมายถึงลักษณะเฉพาะที่เก็บค่าเป็นตัวเลข
  - {V<sub>1</sub>, V, ..., V<sub>n</sub>} หมายถึงลักษณะเฉพาะที่เก็บค่าไม่ต่อเนื่อง
- **@data** เป็นบรรทัดที่บอกถึงแถวที่ตามมาจะเป็นข้อมูล โดยแต่ละแถวจะแทนหนึ่งตัวอย่างข้อมูลซึ่งเรียงตามค่าของลักษณะเฉพาะที่บอกไว้ข้างต้น

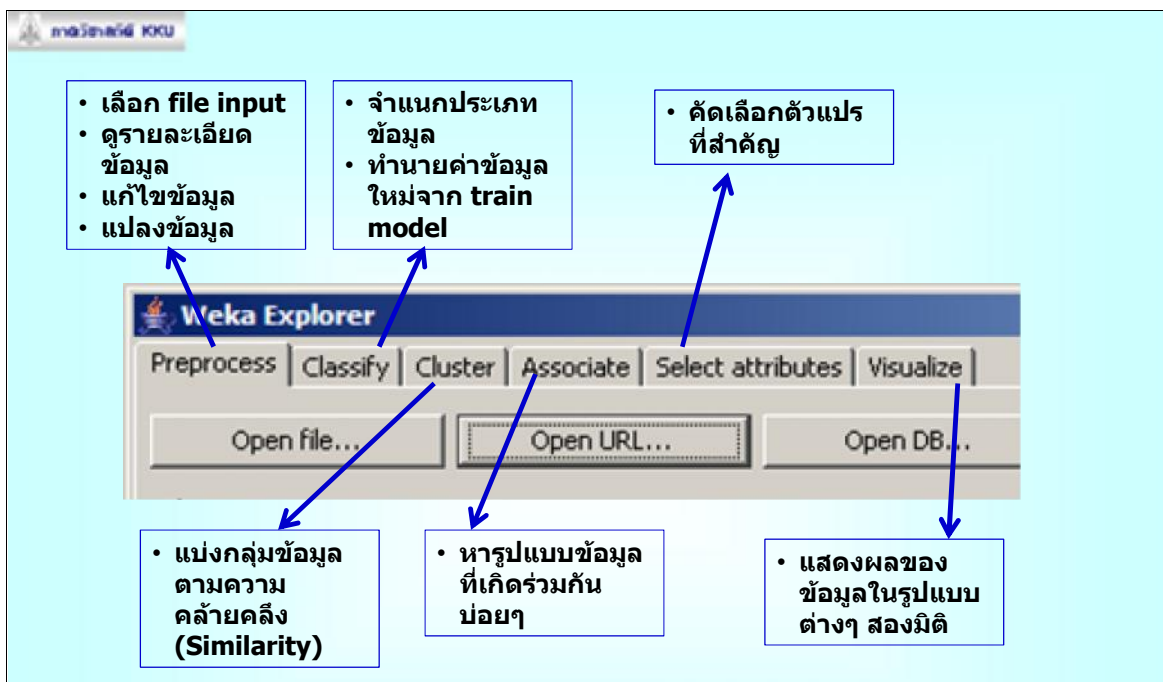
### กรณี CSV file



- สร้างแฟ้มแบบ CSV ด้วยโปรแกรม Microsoft Excel โดยบันทึกสกุลเป็น CSV



## Data Preparation



The screenshot shows the Weka Explorer interface with the 'age' attribute selected. The interface is annotated with yellow boxes and arrows:

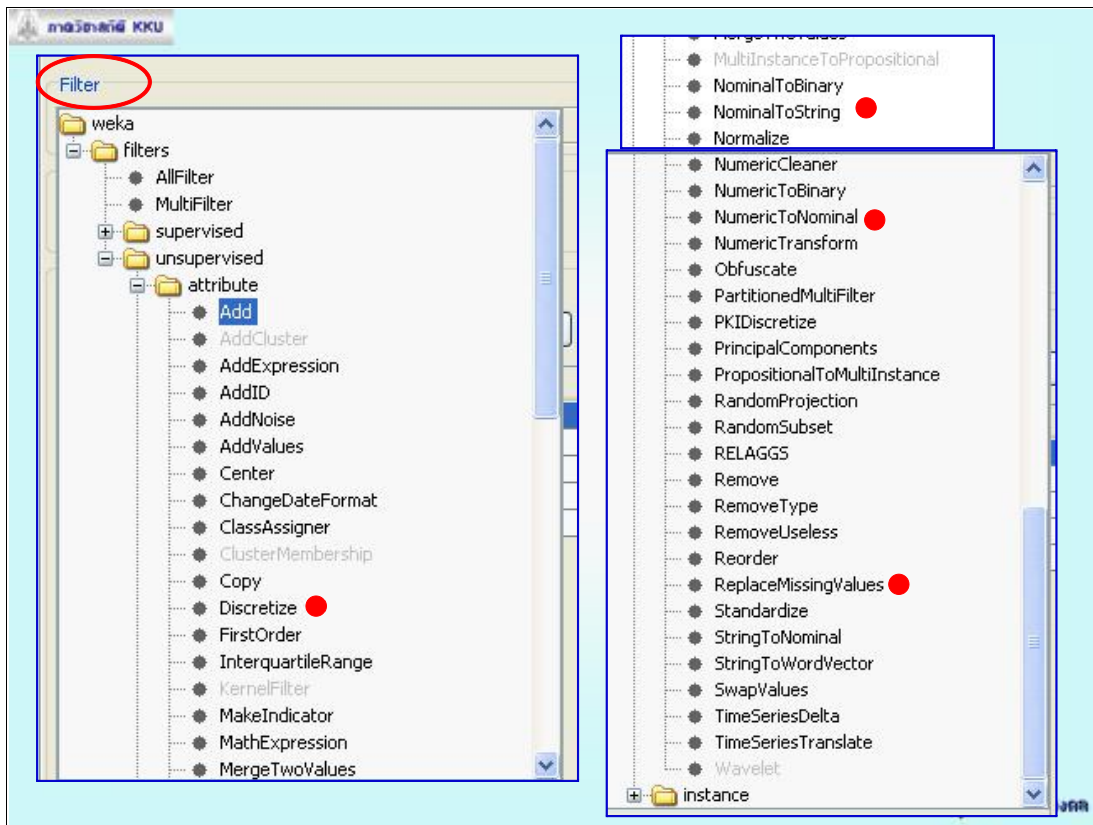
- load:** Points to the 'Open file...', 'Open URL...', and 'Open DB...' buttons.
- filter:** Points to the 'Filter' section, specifically the 'Choose' dropdown set to 'None'.
- analyze:** Points to the 'Attributes' list and the 'Selected attribute' statistics table.

**Selected attribute statistics:**

Statistic	Value
Minimum	17
Maximum	90
Mean	38.452
StdDev	13.598

**Class: class (Nom)**

The histogram shows the distribution of the 'age' attribute, with a red area representing the distribution and a blue area representing the class distribution. The x-axis ranges from 17 to 90, and the y-axis shows frequency.



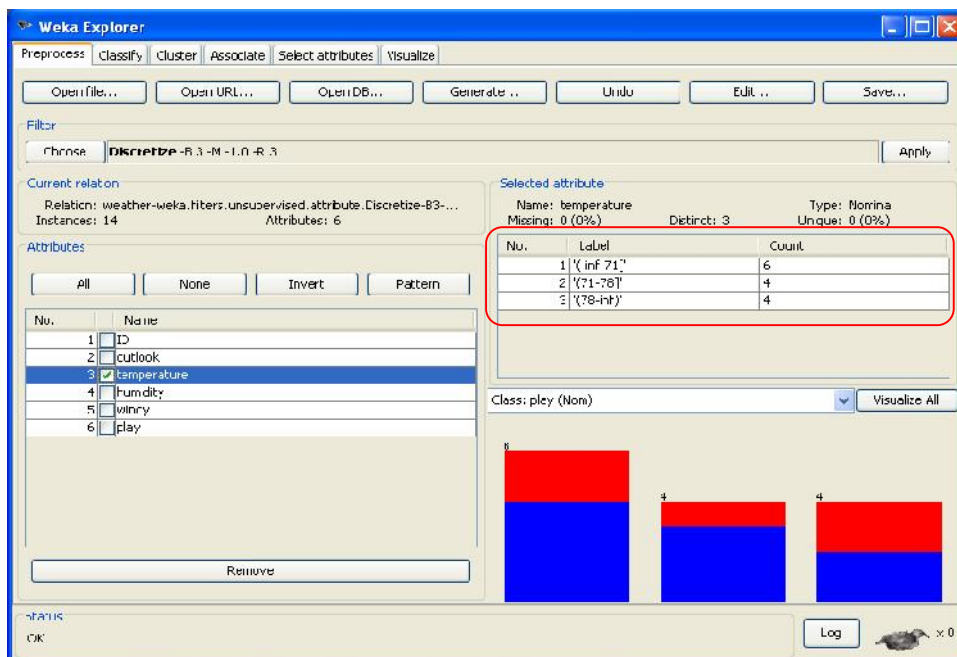
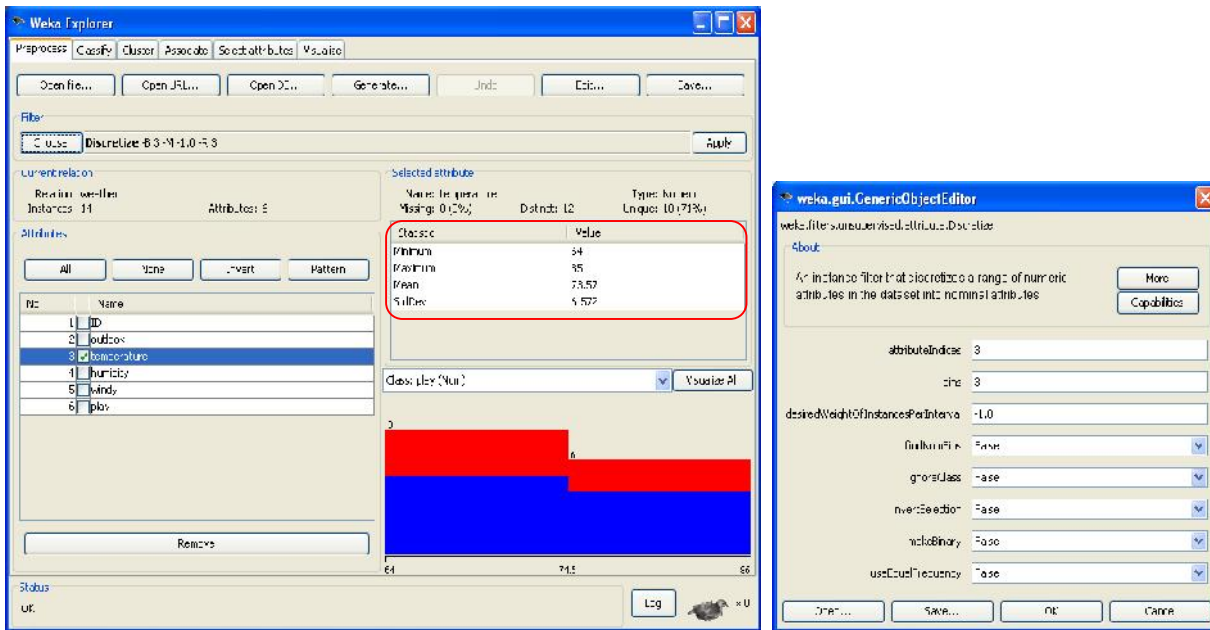
ตัวอย่าง ในการทำ Discretized เป็นการแปลงค่าข้อมูลให้เป็น discrete หรือไม่ต่อเนื่อง เช่น ค่า temperature

### เรื่อง Filter ในส่วนการทำ Discretize

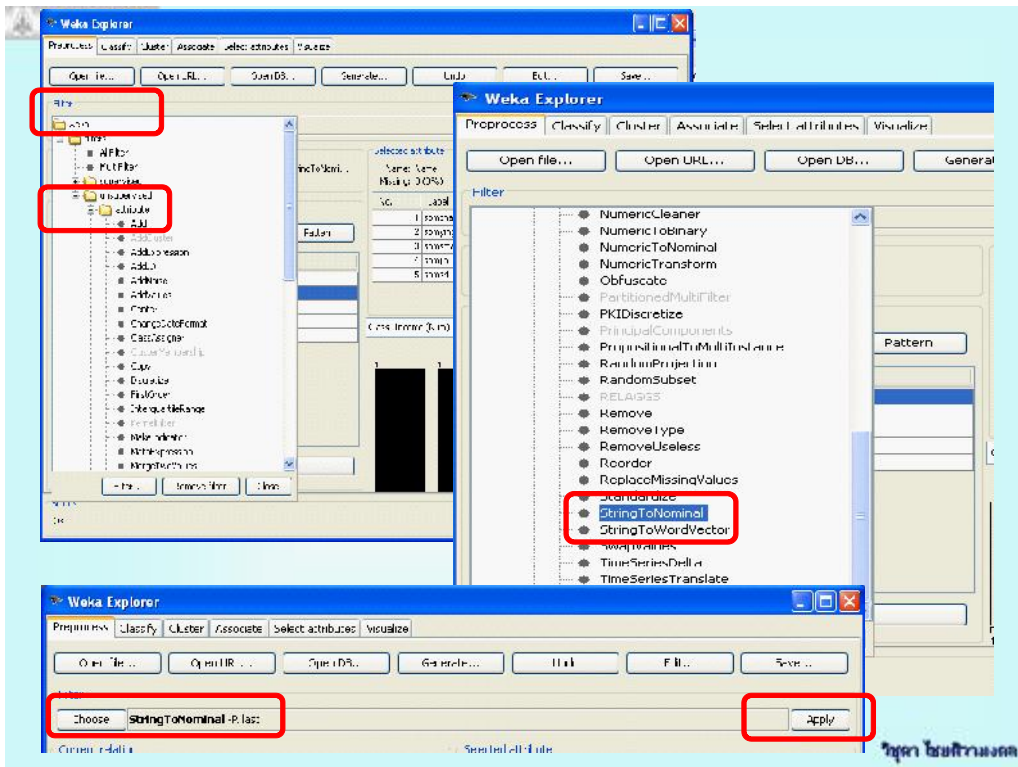
- **Supervised**
  - แปลงข้อมูลแบบอัตโนมัติ
  - ควบคุมด้วยพารามิเตอร์ที่ผู้ใช้กำหนด
- **Unsupervised**
  - แปลงข้อมูลที่ผู้ใช้กำหนดเอง

เป็นการแปลงค่าข้อมูลให้เป็น discrete หรือไม่ต่อเนื่อง โดยผู้ใช้เลือกลำดับของ attribute ที่ต้องการแปลงในช่อง **attribute Indices** และกำหนดจำนวนช่วงที่ต้องการใน **bins** เราสามารถแบ่งแบบ equalwidth หรือ equal depth โดยปรับเป็น False ที่ **useEqualFrequency** หลังจากนั้น กดปุ่ม OK แล้วกดปุ่ม apply

\*สุดา ไชยพิลาภมณฑล







**Supervised**



**1**

No.	outlook	temperature	humidity	windy	play
	Nominal	Numeric	Numeric	Nominal	Nominal
1	sunny	85.0	85.0	FALSE	no
2	sunny	80.0	90.0	TRUE	no
3	overcast	83.0	86.0	FALSE	yes
4	rainy	70.0	96.0	FALSE	yes
5	rainy	68.0	80.0	FALSE	yes
6	rainy	65.0	70.0	TRUE	no
7	overcast	64.0	65.0	TRUE	yes
8	sunny	72.0	95.0	FALSE	no
9	sunny	69.0	70.0	FALSE	yes
10	rainy	75.0	80.0	FALSE	yes
11	sunny	75.0	70.0	TRUE	yes
12	overcast	72.0	65.0	TRUE	yes
13			75.0	FALSE	yes
14			91.0	TRUE	no

**2** click J48

**3** True คือ ให้ตัดแต่งกิ่ง

**4** สังเกตลักษณะข้อมูล เลือก J48 สังเกต Properties -> เลือก Cross-validation Folds 10 -> Click more option- click output prediction

**Bayes → NaïveBayesSimple** **Naïve Bay**

**File 3weather.arff → Class Label คือ Play → ข้อมูลเชิงคุณภาพ 2 level**

No.	outlook Nominal	temperature Numeric	humidity Numeric	windy Nominal	play Nominal
1	sunny	85.0	85.0	FALSE	no
2	sunny	80.0	90.0	TRUE	no
3	overcast	83.0	86.0	FALSE	yes
4	rainy	70.0	96.0	FALSE	yes

**Correctly Classified Instances = 64.2857 %**

**Lazy → IBk** **K-nearest neighbors**

**File 3weather.arff → Class Label คือ Play → ข้อมูลเชิงคุณภาพ 2 level**

No.	outlook Nominal	temperature Numeric	humidity Numeric	windy Nominal	play Nominal
1	sunny	85.0	85.0	FALSE	no
2	sunny	80.0	90.0	TRUE	no
3	overcast	83.0	86.0	FALSE	yes
4	rainy	70.0	96.0	FALSE	yes
5	rainy	68.0	80.0	FALSE	yes
6	rainy	65.0	70.0	TRUE	no

**Correctly Classified Instances = 78.5714 %**

### Functions → MultilayerPerceptron

### Neural Networks

File 3weather.arff → Class Label คือ Play → ข้อมูลเชิงคุณภาพ 2 level

No.	outlook	temperature	humidity	windy	play
	Nominal	Numeric	Numeric	Nominal	Nominal
1	sunny	85.0	85.0	FALSE	no
2	sunny	80.0	90.0	TRUE	no
3	overcast	83.0	86.0	FALSE	yes

**Correctly Classified Instances = 71.4286 %**

### rules → OlexGA

### Genetic Algorithms

File 3weather.arff → Class Label คือ Play → ข้อมูลเชิงคุณภาพ 2 level

การเพิ่ม Module → Olex-GA

**การเพิ่ม module → Olex-GA**  
(a genetic algorithm for the induction of text classification rules)

1. เข้า Web → <https://www.mat.unical.it/Olex-GA/OlexGA/OlexGA-weka.htm>
2. Download file Olex-GA for Weka (binary)
3. เมื่อ extract file จะได้ java 2 file
4. Copy ทั้ง 2 file ลงใน drive C:\Program Files\Weka-3-6
5. เข้าไป edit file configuration ของ weka → ใน RunWeka.ini file เพื่อ include the olexGa.jar and the jaga.jar ลงใน weka classpath ในบรรทัดสุดท้ายเพิ่มคำสั่งดังนี้

`cp=%CLASSPATH%;C:/Program Files/Weka-3-6/jaga.jar;C:/Program Files/Weka-3-6/olexGA.jar`

`cp=%CLASSPATH%;C:/Program Files/Weka-3-6/jaga.jar;C:/Program Files/Weka-3-6/olexGA.jar`

**Functions → LinearRegression** **Linear Regression**

**File CPU → Class ความเร็ว → ข้อมูลเชิงปริมาณ ซึ่งต้องการพยากรณ์**

No.	MYCT	MMIN	MMAX	CACH	CHMIN	CHMAX	class
1	125.0	256.0	6000.0	256.0	16.0	128.0	198.0
2	29.0	8000.0	32000.0	32.0	8.0	32.0	269.0
3	29.0	8000.0	32000.0	32.0	8.0	32.0	220.0
4	29.0	8000.0	32000.0	32.0	8.0	32.0	172.0

Classifier: **LinearRegression -S 0 -R 1.**

Test options:

- Use training set
- Supplied test set
- Cross-validation Folds: **5**
- Percentage split %: **66**

(Num) class: [dropdown]

Start Stop

Result list (right-click for options):

- 21:20:48 - trees.J48
- 21:51:06 - bayes.NaiveBayesSimple
- 22:03:05 - lazy.IBk
- 22:15:07 - functions.MultilayerPerceptron
- 22:25:47 - functions.LinearRegression
- 22:31:10 - functions.LinearRegression

รศ.วิชามงคล

## Unsupervised

**Cluster → SimpleKMeans** **Cluster**

**File... 1bank\_data.arff**

No.	age	sex	region	income	married	children	car	save_act	current_act	mortgage	pep
1	48.0	FEMALE	INNER...	17546.0	NO	1.0	NO	NO	NO	NO	YES
2	40.0	MALE	TOWN	30085.1	YES	3.0	YES	NO	YES	YES	NO
3	51.0	FEMALE	INNER...	16575.4	YES	0.0	YES	YES	YES	NO	NO
4	23.0	FEMALE	TOWN	20375.4	YES	3.0	NO	NO	YES	NO	NO
5	57.0	FEMALE	RURAL	50576.3	YES	0.0	NO	YES	NO	NO	NO
6	57.0	FEMALE	TOWN	37869.6	YES	2.0	NO	YES	YES	NO	YES

Clusterer: **SimpleKMeans -I 6 -A "weka.core.EuclideanDistance"**

Cluster mode:

- Use training set
- Supplied test set
- Percentage split %: **66**
- Classes to clusters evaluation

(Num) pep: [dropdown]

Store clusters for visualization

Ignore attributes: [text box]

Start Stop

weka.gui.GenericObjectEditor:

Cluster data using the k means algorithm.

displayStdDevs:  False

distanceFunction: Choose **EuclideanDistance -R first-last**

dontReplaceMissingValues:  False

maxIterations: 500

**numClusters: 6**

preserveInstancesOrder:  False

seed: 10

Open... Save... OK Cancel

**เทคนิคการวัดความห่าง**

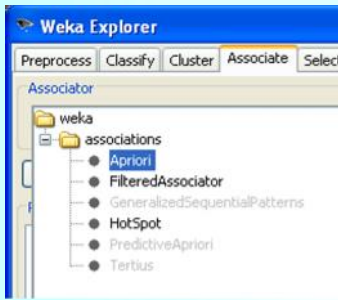
**จำนวน Cluster ที่ต้องการ**

รศ.วิชามงคล

# Associate → Apriori

# Association

File.. 2market\_data.arff

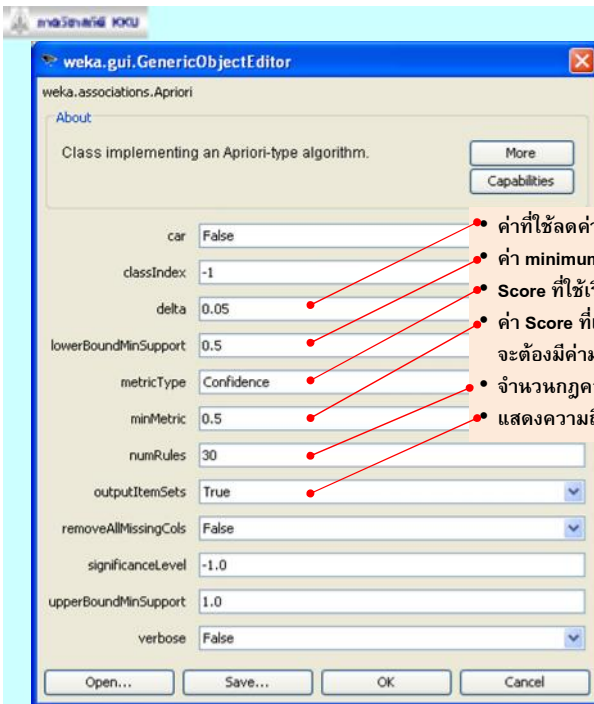


Trans_Id	Items
T1	{Bread, Jelly, PeanutButter}
T2	{Bread, PeanutButter}
T3	{Bread, Milk, PeanutButter}
T4	{Beer, Bread}
T5	{Beer, Milk}



No.	ID	Bread	Jelly	PeanutButter	Milk	Beer
1	T1	y	y	y		
2	T2	y		y		
3	T3	y		y	y	
4	T4	y				y
5	T5				y	y

รศ.วิชาญมงคล



- ค่าที่ใช้ลดค่า minimum support
- ค่า minimum support
- Score ที่ใช้เรียง (rank) กฎความสัมพันธ์
- ค่า Score ที่เลือกใน metric Type โดยกฎที่สนใจจะต้องมีค่ามากกว่าที่กำหนด
- จำนวนกฎความสัมพันธ์ที่ต้องการ
- แสดงความถี่ของสินค้าที่ต้องการซื้อพร้อมกัน

รศ.วิชาญมงคล

# Associate → GeneralizedSequentialPatterns

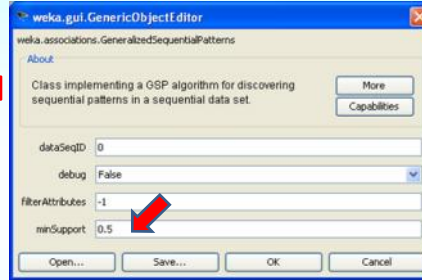
File Asso\_Sequential2.csv

	A	B	C
1	ID	time	type
2	R1	1	fruit
3	R1	2	cloth
4	R2	1	shoes
5	R2	2	cloth
6	R2	3	fruit
7	R3	1	pen
8	R3	2	cloth
9	R4	1	shoes
10	R4	2	cloth
11	R5	1	cloth
12	R5	2	fruit
13	R5	3	pen
14			

Viewer

Relation: Asso\_Seq1-weka.filters.u

No.	ID	time	type
1	R1	1	fruit
2	R1	2	cloth
3	R2	1	shoes
4	R2	2	cloth
5	R2	3	fruit
6	R3	1	pen
7	R3	2	cloth
8	R4	1	shoes
9	R4	2	cloth
10	R5	1	cloth
11	R5	2	fruit
12	R5	3	pen



ต้องประกอบด้วย ID และ Time\_Stamp

วิศวะ ไซวิทยา มงคล

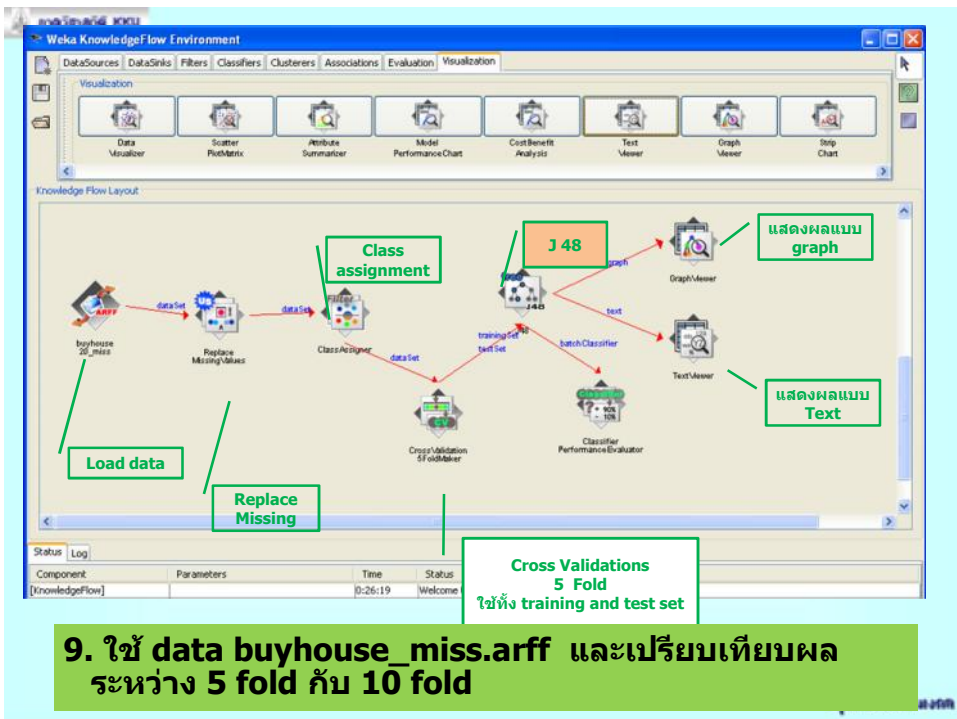
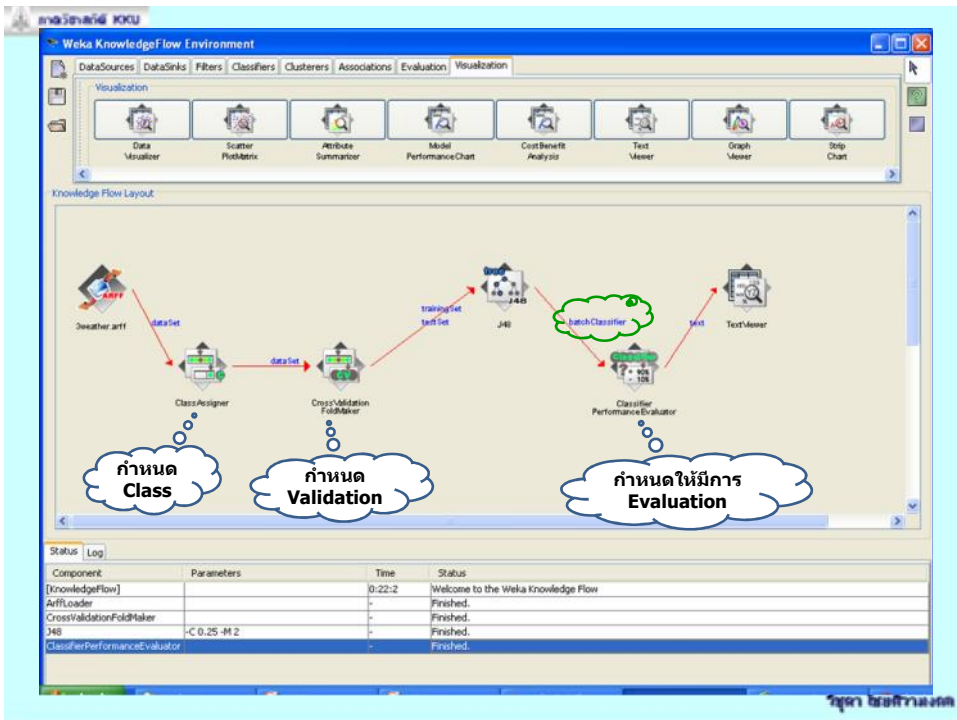
## KnowledgeFlow

Preprocess

ใช้ประเมินตัวแบบ

ใช้แสดงผล

วิศวะ ไซวิทยา มงคล



**9. ใช้ data buyhouse\_miss.arff และเปรียบเทียบผล ระหว่าง 5 fold กับ 10 fold**