# Slide 1

# Data Mining

*(Knowledge discovery in database)*

# Slide 2

# What is Data Mining?

- *Data Mining: "The non trivial extraction of implicit, previously unknown, and potentially useful information from data"*
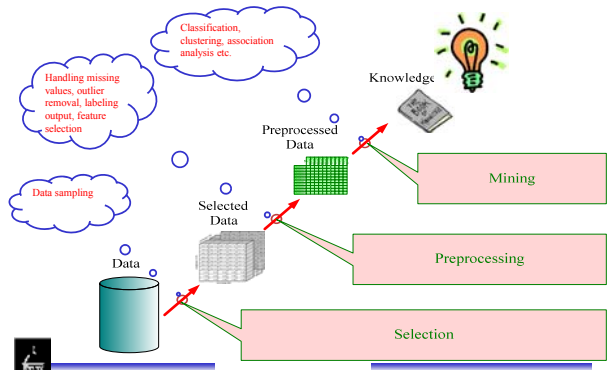  - *William J Frawley, Gregory Piatetsky-Shapiro and Christopher J Matheus*

- **Data mining finds valuable information hidden in large volumes of data.**

- **Data Mining process involves:**
  - Databases
  - Statistics
  - Machine Learning
  - High Performance Computing
  - Visualization
  - Mathematics

# Slide 3

# Data mining: Basic steps

Classification, clustering, association analysis etc.

Handling missing values, outlier removal, labeling output, feature selection

Data sampling

Knowledge

Preprocessed Data

Mining

Selected Data

Preprocessing

Data

Selection

# Slide 4

# Mining tasks

- **Classification**: *YES, NO*    Patient diagnosis

- **Regression**: *Predict the actual output value*    Weather forecasting

- **Clustering**: '*Grouping identical items*'    Recommender systems e.g. Netflix

- **Association rules**: *Identifying association among input attributes.*    {Milk} --> {Coke}, {Diaper, Milk} --> {Beer}

- **Anomaly detection**: *Detect deviation from normal behavior*    Fraud detection, Network intrusions

## Data Mining Software

- **Enterprise-level: (US $10,000 and more)**
  - Fair Isaac, IBM, Insightful, KXEN, Oracle, SAS, and SPSS
- **Department-level: (from $1,000 to $9,999)**
  - Angoss, CART/MARS/TreeNet/Random Forests, Equbits, GhostMiner, Gornik, Mineset, MATLAB, Megaputer, Microsoft SQL Server, Statsoft Statistica, ThinkAnalytics
- **Personal-level: (from $1 to $999):**
  - Excel, See5, MATLAB
- **Free:**
  - C4.5, R, **Weka**, Xelopes

## Data Mining: WEKA

## Outline

- Data preparation
- Preprocessing and "arff" files
- Filters, classifiers, and visualization
- Attribute selection
- Training and testing
- Quality measurements
- Interpretation of results

## Data Mining Procedures

- Prepare the data into desired formats
- Preprocess the data if necessary
- Select different algorithms based on application or domain expertise
- Evaluate the results and repeat experiments again if necessary

## Data Sets



The University of Iowa     Intelligent Systems Laboratory

## Arff file

### Header of arff file

| @relation | weather | |
|---|---|---|
| @attribute | No | real |
| @attribute | outlook | {sunny,overcast,rainy} |
| @attribute | temperature | real |
| @attribute | humidity | real |
| @attribute | windy | {TRUE,FALSE} |
| @attribute | play | {yes,no} |

The University of Iowa     Intelligent Systems Laboratory

## Arff file

### Data in arff file

Weather_training - Not...

File   Edit   Format   View   Help

```
@data
1,sunny,85,85,FALSE,no
2,sunny,80,90,TRUE,no
3,overcast,83,86,FALSE,yes
4,rainy,70,96,FALSE,yes
5,rainy,68,80,FALSE,yes
6,rainy,65,70,TRUE,no
7,overcast,64,65,TRUE,yes
8,sunny,72,95,FALSE,no
9,sunny,69,70,FALSE,yes
10,rainy,75,80,FALSE,yes
11,sunny,75,70,TRUE,yes
12,overcast,72,90,TRUE,yes
13,overcast,81,75,FALSE,yes
14,rainy,71,91,TRUE,no
```

The University of Iowa     Intelligent Systems Laboratory

## Arff file



### Training data set

The University of Iowa     Intelligent Systems Laboratory

3

## WEKA Explorer

## WEKA

## Parameter Distribution

## Filters

# Filters



The University of Iowa                    Intelligent Systems Laboratory

# Filters



The University of Iowa                    Intelligent Systems Laboratory
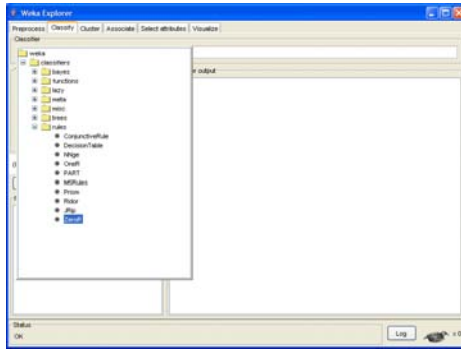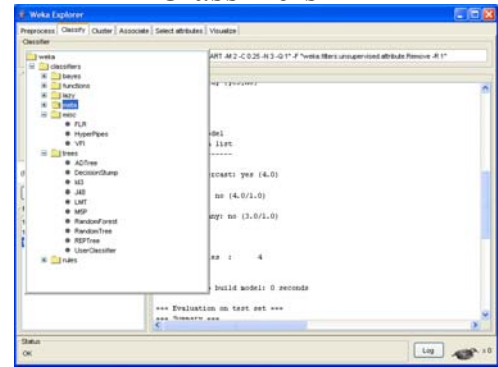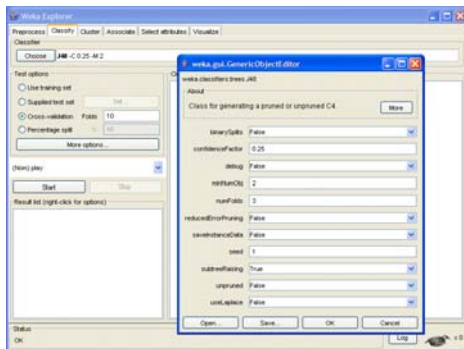
# Filters



The University of Iowa                    Intelligent Systems Laboratory

# Filters



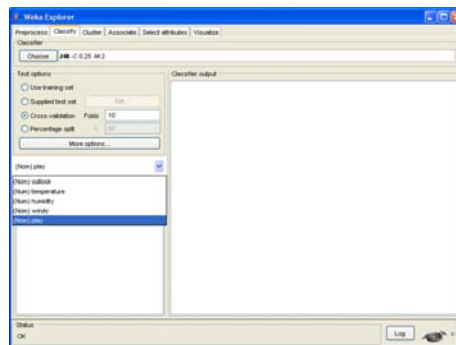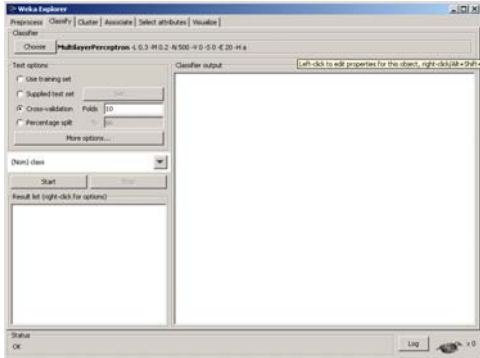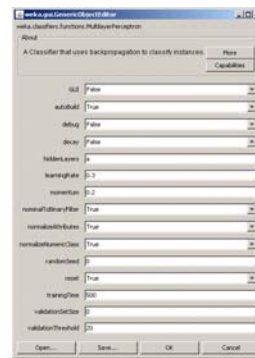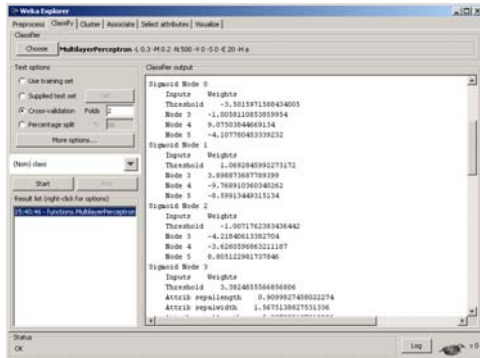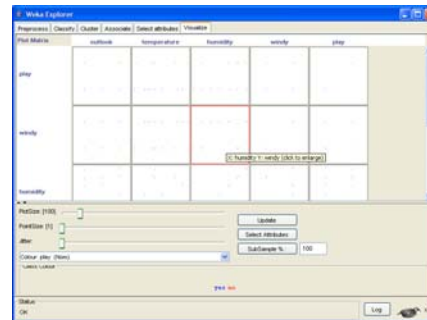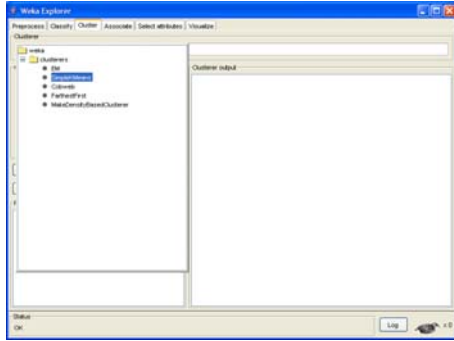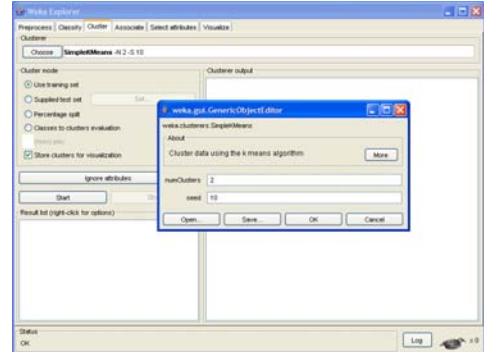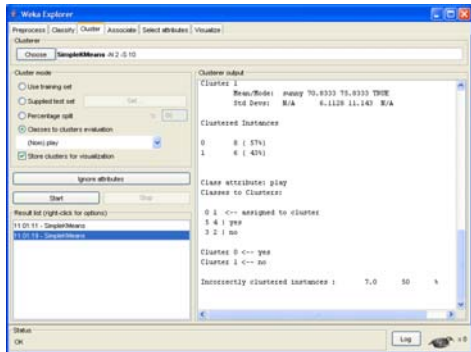The University of Iowa                    Intelligent Systems Laboratory

## Classifier

## Classifiers

## Decision Tree

## Decision Tree

# Decision Tree



# Decision Tree



# Decision Tree



# Decision Tree

# Neural Networks



# Neural Networks



# Neural Networks



# Visualization

# Clustering

# Clustering

# Clustering

# UCI repository

- http://archive.ics.uci.edu/ml/index.html