# An Introduction to WEKA 3.9.x

Some of the slides are taken from presentation by Yizhou Sun

1

1

# Content

- What is WEKA?
- The Explorer:
  - Preprocess data
  - Classification
  - Clustering
  - Association Rules
  - Attribute Selection
  - Data Visualization
- References and Resources

2

2

# What is WEKA?

- **W**aikato **E**nvironment for **K**nowledge **A**nalysis
  - It's a data mining/machine learning tool developed by Department of Computer Science, University of Waikato, New Zealand.
  - Weka is also a bird found only on the islands of New Zealand.
  - https://www.youtube.com/watch?v=1vgA3CN2PH0
  - Weka software is developed in Java.

3

3

# Download and Install WEKA

- Website: http://www.cs.waikato.ac.nz/~ml/weka/index.html
- Support multiple platforms (written in java):
  - Windows, Mac OS X and Linux

4

4

## Main Features

- 49+ data preprocessing tools
- 76+ classification/regression algorithms
- 8+ clustering algorithms
- 3+ algorithms for finding association rules
- 15+ attribute/subset evaluators + 10+ search algorithms for feature selection

5

5

## Main GUI

- Four graphical user interfaces
  - "The Explorer" (exploratory data analysis)
  - "The Experimenter" (experimental environment)
  - "The KnowledgeFlow" (new process model inspired interface)
  - "Workbench" (unified GUI that combines above three)
- One old fashioned Command Line Interface (CLI)



6

6

## The package management system

- Weka community keeps adding new algorithms and features.
- These are placed into plugin packages.
- A package management system allows the user to browse and install packages of interest.



7

7

## Content

- What is WEKA?
- The Explorer:
  - Preprocess data
  - Classification
  - Clustering
  - Association Rules
  - Attribute Selection
  - Data Visualization
- References and Resources

8

8

## Explorer: pre-processing the data

- Data can be imported from a file in various formats: ARFF, CSV, C4.5, binary
- Data can also be read from a URL or from an SQL database (using JDBC)
- Pre-processing tools in WEKA are called "filters"
- WEKA contains filters for:
  - Discretization, normalization, resampling, attribute selection, transforming and combining attributes, …

9

9

## Weka explorer interface



10

10

## WEKA: Attribute Relation File Format (ARFF)

CSV format

ARFF format



11

11

## Load file weather.numeric.arff



12

12

## Building a decision tree – Select J48 (implementation of C4.5 algo)



13

## Using C4.5 classifier



14

## C4.5 Decision tree visualizer



15

15

## Using training set for testing



16

16

17



18

## Working with iris data – iris.arff



19

## Use J48 with cross validation



20

## Iris classifier visualization

Weka Classifier Visualize: 15:29:22 - trees.J48 (iris)

X: sepallength (Num)    Y: sepalwidth (Num)

Colour: class (Nom)    Select Instance

Reset   Clear   Open   Save    Jitter

Plot: iris_predicted

Weka: Insta...

Plot : weka.classifiers.trees.J48 (iris)
Instance: 92
        sepallength : 4.9
         sepalwidth : 2.5
        petallength : 4.5
         petalwidth : 1.7
prediction margin : -1.0
   predicted class : Iris-versicolor
             class : Iris-virginica

Class colour

Iris-setosa        Iris-versicolor        Iris-virgin

21

21

## Iris data – J48 tree

Weka Classifier Tree Visualizer: 15:29:22 - trees.J48 (iris)

Tree View

petalwidth

<= 0.6          > 0.6

Iris-setosa (50.0)          petalwidth

<= 1.7          > 1.7

petallength          Iris-virginica (46.0/1.0)

<= 4.9          > 4.9

Iris-versicolor (48.0/1.0)          petalwidth

<= 1.5          > 1.5

Iris-virginica (3.0)          Iris-versicolor (3.0/1.0)

22

22

# Open CSV file in Weka



23

23

# Filter example – Remove attribute



24

24

## Iris data with first attribute removed



25

## Load CPU data – cpu.arff



26

# Choose M5P Quinlan's M5 algorithm
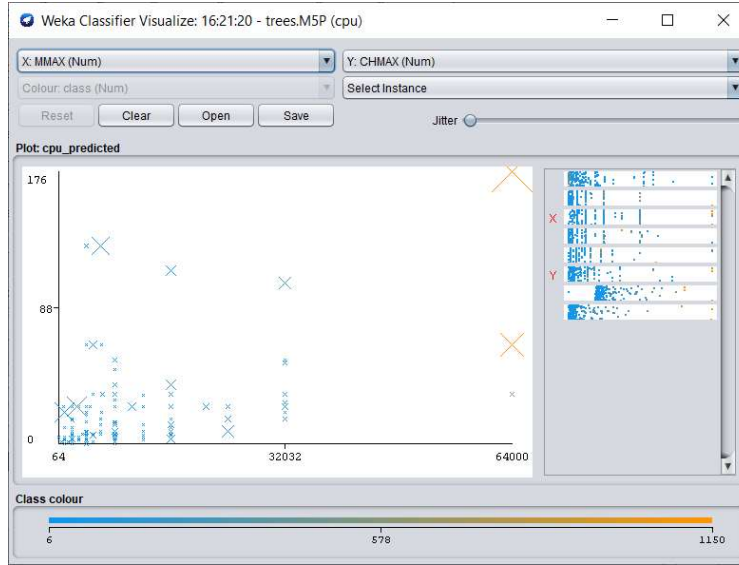# Decision trees of regression models



27

# M5P linear regression models tree



- The first number at each leaf is the number of instances that reach it
- The second is the root mean squared error of the predictions expressed as a percentage of the standard deviation

28

## Visualization of errors – larger the cross, larger the error – M5P



29

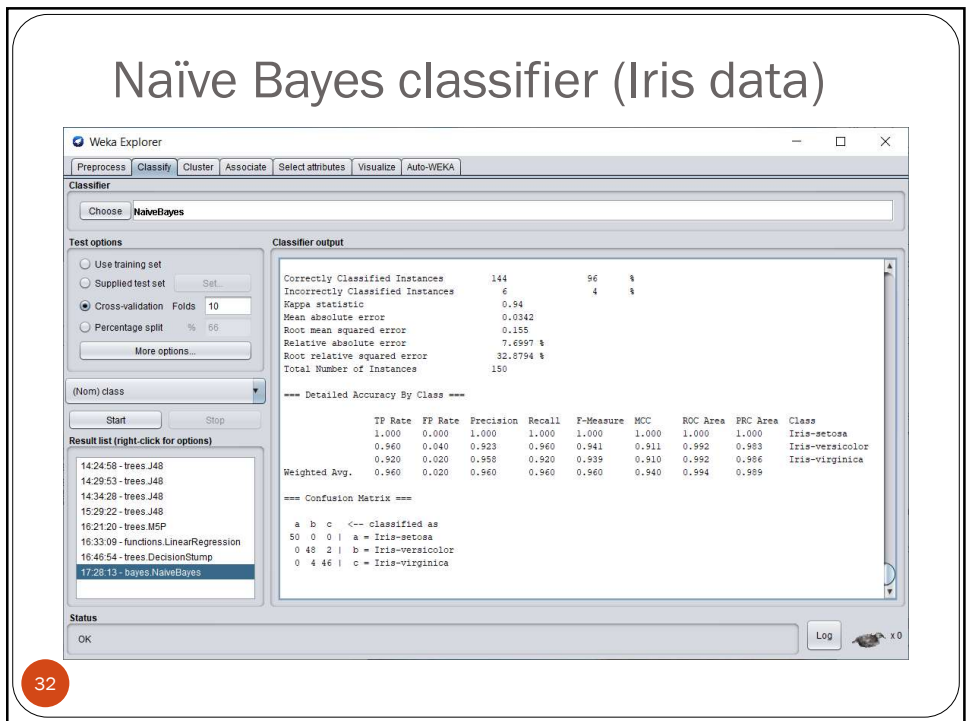## Choose functions -> Linear regression for CPU data



30

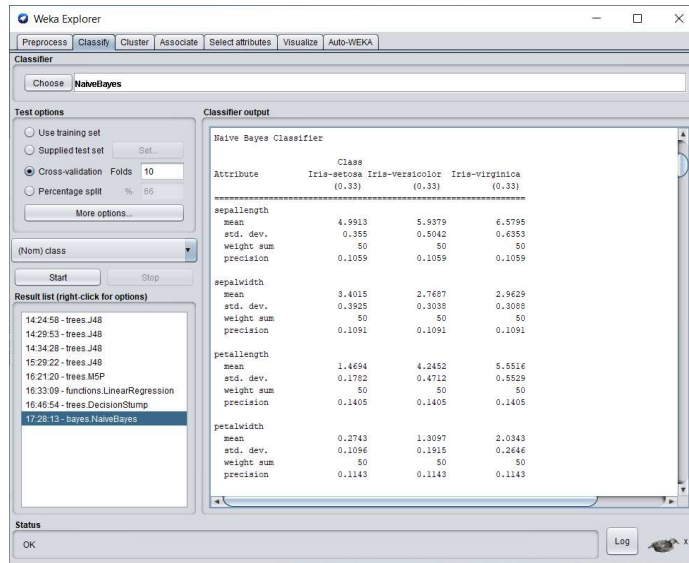# Visualization of errors – linear regression (M5P is superior)


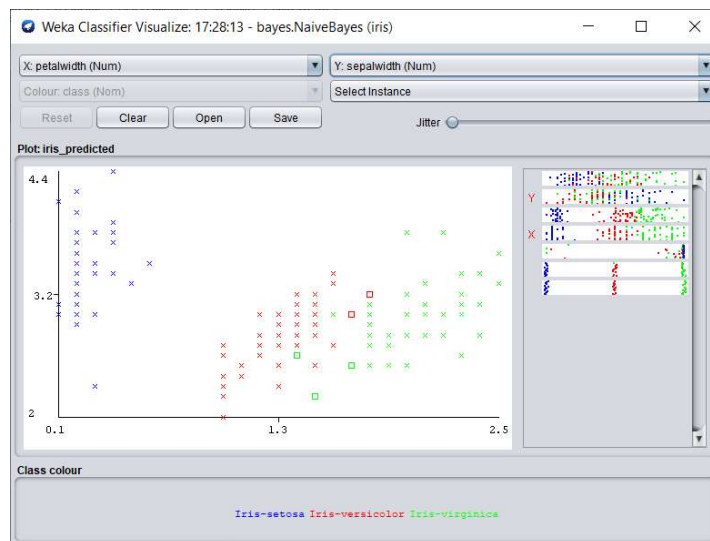
31

31

# Naïve Bayes classifier (Iris data)



32

32

## Naïve Bayes classifier uses normal distribution to model numeric attributes



33

## Naïve Bayes classifier visualization



34

## Naïve Bayes updatable – process one instance at a time



35

## Naïve Bayes updatable visualization



36

## Explorer: clustering data

- WEKA contains "clusterers" for finding groups of similar instances in a dataset
- Implemented schemes are:
  - *k*-Means, EM, Cobweb, *X*-means, FarthestFirst
- Clusters can be visualized and compared to "true" clusters (if given)
- Evaluation based on loglikelihood if clustering scheme produces a probability distribution

37

37

## The K-Means Clustering Method

- Given *k*, the *k-means* algorithm is implemented in four steps:
  - Partition objects into *k* nonempty subsets
  - Compute seed points as the centroids of the clusters of the current partition (the centroid is the center, i.e., *mean point*, of the cluster)
  - Assign each object to the cluster with the nearest seed point
  - Go back to Step 2, stop when no more new assignment

38

38

## Clustering – open iris.arff, select Simple Kmeans (numClusters = 3)



39

## K-means clustering visualization



40

## EM (Expectation-Maximization) clustering – out of the box



41

41

## EM clustering (out of the box) visualization



42

42

EM clustering – numClusters = 3

43



EM (numClusters = 3) visualization

44

## Hierarchical clustering (numClusters = 3)



45

## Hierarchical clustering – tree visualization



46

## Hierarchical clustering - visualization



47

## Explorer: finding associations

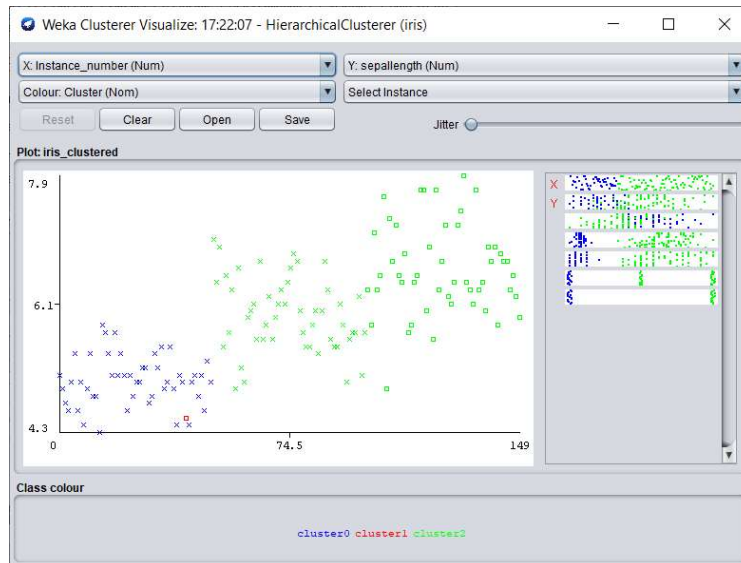- WEKA contains an implementation of the Apriori algorithm for learning association rules
  - Works only with discrete data
- Can identify statistical dependencies between groups of attributes:
  - milk, butter $\Rightarrow$ bread, eggs (with confidence 0.9 and support 2000)
- Apriori can compute all rules that have a given minimum support and exceed a given confidence
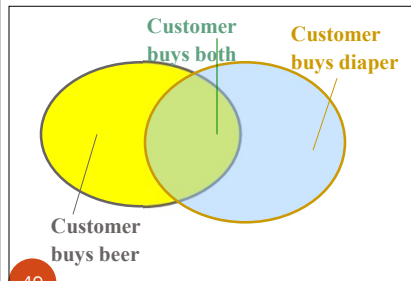
48

## Basic Concepts: Frequent Patterns

| Tid | Items bought |
|-----|--------------|
| 10 | Beer, Nuts, Diaper |
| 20 | Beer, Coffee, Diaper |
| 30 | Beer, Diaper, Eggs |
| 40 | Nuts, Eggs, Milk |
| 50 | Nuts, Coffee, Diaper, Eggs, Milk |

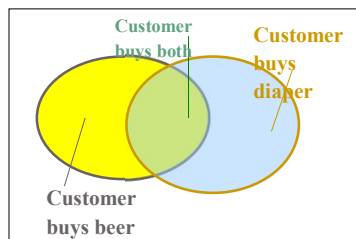**Customer buys both**

**Customer buys diaper**

**Customer buys beer**

49

- itemset: A set of one or more items
- k-itemset $X = \{x_1, \ldots, x_k\}$
- *(absolute) support*, or, *support count* of X: Frequency or occurrence of an itemset X
- *(relative) support*, *s*, is the fraction of transactions that contains X (i.e., the probability that a transaction contains X)
- An itemset X is *frequent* if X's support is no less than a *minsup* threshold

49

## Basic Concepts: Association Rules

| Tid | Items bought |
|-----|--------------|
| 10 | Beer, Nuts, Diaper |
| 20 | Beer, Coffee, Diaper |
| 30 | Beer, Diaper, Eggs |
| 40 | Nuts, Eggs, Milk |
| 50 | Nuts, Coffee, Diaper, Eggs, Milk |

**Customer buys both**

**Customer buys diaper**

**Customer buys beer**

50

- Find all the rules $X \rightarrow Y$ with minimum support and confidence
  - support, *s*, probability that a transaction contains $X \cup Y$
  - confidence, *c*, conditional probability that a transaction having X also contains Y

*Let minsup = 50%, minconf = 50%*

*Freq. Pat.:* Beer:3, Nuts:3, Diaper:4, Eggs:3, {Beer, Diaper}:3

- Association rules: (many more!)
  - *Beer* → *Diaper* (60%, 100%)
  - *Diaper* → *Beer* (60%, 75%)

50

## Associative rules: open supermarket.arff



51

## Associative rules: Apriori algorithm



52

## Associative rules – Apriori output



53

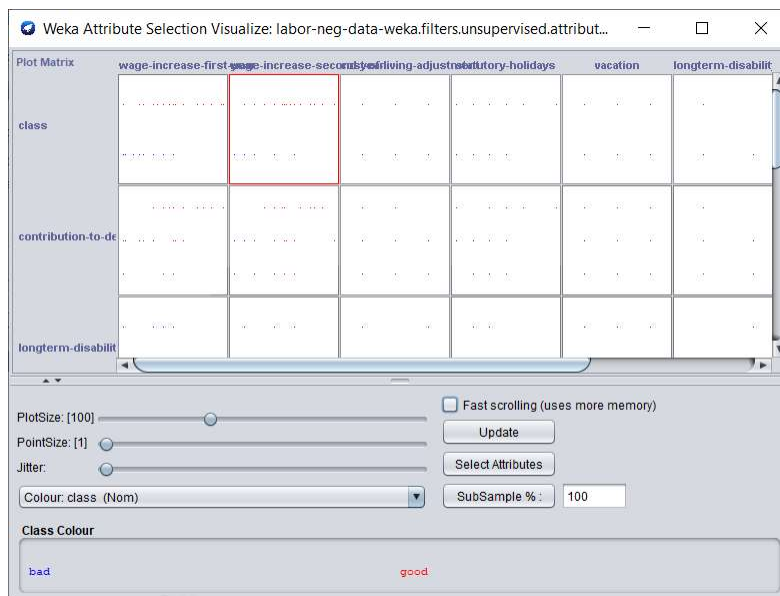## Weka – Automated feature selection – labor.arff



54

## Weka – Automated feature selection
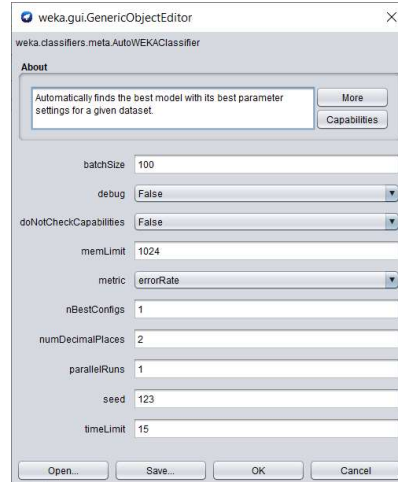## Select attributes

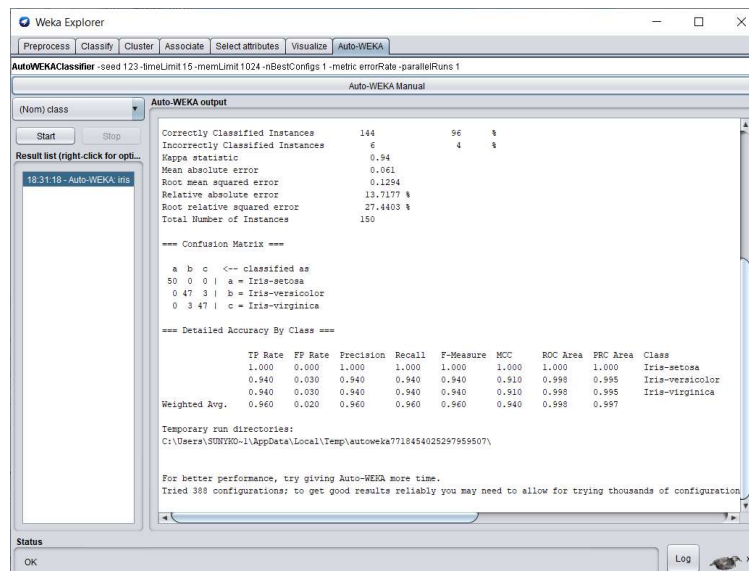

55

## Attribute selection visualize



56

## AutoWeka

- Install AutoWeka plugin from Weka plugin manager
- Automatically searches through the joint space of WEKA's learning algorithms and their respective hyperparameter settings to maximize performance
- Load iris.arff
- Tune AutoWeka parameters
- Hit start

57

57

## AutoWeka result

58

58

## Imbalanced class problem

- Use Weka SMOTE filter
- Example of German Credit Data Analysis
- Two types of risks are associated with the bank's decision:
- If the applicant is a good credit risk, i.e. is likely to repay the loan, then not approving the loan to the person results in a loss of business to the bank
- If the applicant is a bad credit risk, i.e. is not likely to repay the loan, then approving the loan to the person results in a financial loss to the bank
- The German Credit data set is a publically available data set downloaded from the UCI Machine Learning Repository

59

59

## Imbalanced class problem – German_credit dataset



60

60

# Use classbalancer from filter and apply



61

# Use classbalancer from filter and apply



62

# Dealing with missing values – labor negotiation dataset



63

63

# Filters -> unsupervised -> attribute -> replaceMissingValues -> Apply



64

64

# Replace numeric values with mean and nominal values with mode



65

65

# Try Naïve Bayes on modified data



66

66

# References and Resources

- References:
  - WEKA website: http://www.cs.waikato.ac.nz/~ml/weka/index.html
  - WEKA Tutorial:
    - Machine Learning with WEKA: A presentation demonstrating all graphical user interfaces (GUI) in Weka.
    - A presentation which explains how to use Weka for exploratory data mining.
  - WEKA Data Mining Book:
    - Ian H. Witten and Eibe Frank, Data Mining: Practical Machine Learning Tools and Techniques (Fourth Edition)
  - WEKA Wiki: http://weka.sourceforge.net/wiki/index.php/Main_Page
  - AutoWeka Software:
    http://www.cs.ubc.ca/labs/beta/Projects/autoweka/#software
  - Others:
    - Jiawei Han and Micheline Kamber, Data Mining: Concepts and Techniques, 2nd ed.

67