# Data Mining

Module:
Data Mining with Weka

Last updated 9/20/2020

# Overview of DM with Weka Module

- Weka-1: Weka Basics
- Weka-2: Preprocessing
- Weka-3: Decision Trees
- Weka-4: Other Classifiers & Feature Selection
- Weka-5: KnowledgeFlow Interface

# Weka-1: Weka Basics

- What is Weka?
- Why use Weka (and why not)
- Weka versions and downloading the software
- Weka ARFF file format
- Starting Weka and the 5 Weka Interfaces

# What is Weka?

- Weka is a data mining suite developed at University of Waikato
- Weka stands for Waikato Environment for Knowledge Analysis
- Weka includes everything necessary to generate and apply data mining models
  - Covers all major data mining tasks
  - Includes tools to preprocess and visualize data
  - Includes multiple (5) interfaces
    - We will focus on the explorer interface
    - Briefly discuss the knowledgeflow
    - Does not require any programming

# WEKA is also a bird



*Copyright: Martin Kramer (mkramer@wxs.nl)*

# Why Weka?

- There are many options on what data mining suite or toolkit one can use
- Weka has the following advantages:
  - Easy to learn
  - Free and open-source
  - Easy to download and runs on many platforms
  - Does not require programming knowledge
    - This is important since a few students may not have much programming experience

# Weka Disadvantages

- Weka is not as flexible as other tools that are programming based
  - If you are programming anyway you can combine capabilities more flexibly and write code to do things not supported by Weka
  - Tools based on Python or R have large ecosystems

# DM Tools for This Course

- For your project you may use Weka, Python, or even another data mining toolkit

- I do provide snippets of Python code and Weka examples in some lectures
  - But everyone should learn WEKA and complete the tutorial/exercises in this module

- We will not cover the details of how to use Python in this course

# WEKA: the software

- Data mining software written in Java
  - distributed under the GNU Public License
- Used for research, education, and applications
- Main features:
  - Comprehensive set of data pre-processing tools, learning algorithms and evaluation methods
  - Multiple interfaces that do not require programming
  - Experimenter can compare learning algorithms

# Weka: Versions

- There are several versions of WEKA:
  - As of Feb 2020 the stable version is 3.8.4 and that is the one you should be using
- Many of these slides are based on older version and will look different from what you see
  - Dr. Weiss has added notes for significant differences
- Class WEKA page provides relevant info and links
  - https://storm.cis.fordham.edu/~gweiss/data-mining/weka.html
- Can also use following link to download software and documentation
  - https://www.cs.waikato.ac.nz/ml/weka/

# Downloading Weka

- Weka runs on all major platforms
- It is free and easy to download
- Download Weka from here:
  - [https://waikato.github.io/weka-wiki/downloading_weka/](https://waikato.github.io/weka-wiki/downloading_weka/)
  - Use the latest stable version (3.8.4 as of Aug 2020)
- There is plenty of documentation
  - [https://waikato.github.io/weka-wiki/documentation/](https://waikato.github.io/weka-wiki/documentation/)

# Weka ARFF File Format

- ARFF= Attribute Relation File Format
- Weka ARFF files include two main parts:
    - Specification of the features
    - The actual data
- Files are "flat files" usually comma separated
- Other tools use a separate file for each part
    - C4.5 decision tree tool, which was once very popular, had a "names" and "data" file
    - I find the single file format odd, since specification is short but actual data can be millions of lines

# Weka ARFF Example

This just defines dataset name

@relation heart-disease-simplified

Categorical feature must list values

@attribute age numeric
@attribute sex { female, male}
@attribute chest_pain_type { typ_angina, asympt, non_anginal, atyp_angina}
@attribute cholesterol numeric
@attribute exercise_induced_angina { no, yes}
@attribute class { present, not_present}

@data
63,male,typ_angina,233,no,not_present
67,male,asympt,286,yes,present
67,male,asympt,229,yes,present
38,female,non_anginal,?,no,not_present
...

# Starting Weka

- Run Weka 3.8.4 (with console)
  - Different version if a new stable version released
  - May want to create a shortcut on your desktop



This is the GUI Chooser Window

# The Weka Interfaces

- Explorer:
    - We cover in detail
    - Enables you to apply multiple actions but does not explicitly model them
- Experimenter
    - Run many experiments in controlled manner and compare results
- KnowledgeFlow:
    - Like Explorer but each step is represented as a node in a graph, so flow explicitly represented
- Workbench
    - Combines all GUI interfaces into one
- Simple CLI (command line interface)
    - No GUI so can use shell scripts to control experiments
    - Similar to using Python where each command is a function

# Weka-1 Review: Weka Basics

- You should know the following:
  - What Weka is
  - The benefits & drawbacks of WEKA
  - How to download WEKA
    - Ideally you should have installed it
  - The ARFF format and understand it
  - The 5 WEKA interfaces and what they are used for

# Weka-2: Preprocessing

- How to import data
- Visualizing feature value distribution and relationships to class values
- Preprocessing filters
  - Example using "Discretize" filter
- It is strongly recommended that you follow along on your own installation of Weka
  - Reproduce each step on your computer

# Start the Explorer

Click on "Explorer " Button in the GUI Chooser Window

# Importing Data

- Data can be imported from various formats:
  - ARFF, CSV, C4.5
  - You will most often import data in csv if not already prepared for WEKA
- Weka steps
  - Make sure Explorer "PreProcess" tab is selected
  - Click "Open File…" or "Open URL…"
  - Set the extension appropriately (default is .arff) and navigate to the file and select it
- If first line has feature names they are used
  - If no feature names then generic names are generated
  - Add feature names if missing or models won't be interpretable

# Reading in the Iris Dataset

- We start with the well-known Iris data set
  - Weka probably installed it so search for "iris".
  - Otherwise download it from my copy on internet:
    - https://storm.cis.fordham.edu/~gweiss/data-mining/weka-data/iris.arff

- Open it using "Open File" or "Open URL"

- Weka immediately shows stats about the features and how they are distributed

- Other data sets it you want to play around:
    - https://storm.cis.fordham.edu/~gweiss/data-mining/datasets.html
  - Hundreds of others available from UCI Machine Learning Repository
    - https://archive.ics.uci.edu/ml/

Distribution of the 3 classes, which are evenly distribute

Visualization of all feature values with class information. This makes it clear that most features are useful for distinguishing classes.

# Preprocessing Filters

- Pre-processing tools in WEKA called "filters"
- Many preprocessing filters available:
  - Discretization
  - Normalization
  - Resampling
  - Feature selection
  - Feature transformation
  - Many more
- Know what available so can use it when needed
  - We will focus on Discretization

Clicking on this will bring up a list of all of the filters, organized into a hierarchy. Click on each folder to expand the list. There are dozens of choices

# Discretization

- Discretization: numerical feature → categorical
  - In preprocess tab select "Choose", expand "unsupervised" then "attribute"
    - Select "Discretize" filter
    - Note discretize applies to attributes not instances
    - The following few slides are from an earlier version and may look just a little bit different
    - Note: there is a version under the "supervised" folder but it works differently and has different options

**Create 10 equal frequency bins**
Will do for all numerical features given "attributeIndices" value

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Undo | Save...

Filter
Choose | Discretize -B 10 -R first-last | Apply

Current relation
Relation: iris
Instances: 150        Attributes: 5

weka.gui.GenericObjectEditor
weka.filters.unsupervised.attribute.Discretize

About
An instance filter that discretizes a range of numeric attributes in the dataset into nominal attributes.    More

Attributes

No.        Name
1 sepallength
2 sepalwidth
3 petallength
4 petalwidth
5 class

attributeIndices   first-last
bins   10
findNumBins   False
invertSelection   False
makeBinary   False
useEqualFrequency   False

Click for documentation

1) Toggle to "True"

Open...   Save...   OK   Cancel

1        3.95        6.9

Status
OK        Log   x 0

2) Click OK

Weka Knowledge Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Undo | Save...

**Filter**

Choose | **Discretize –F –B 10 –R first–last** | Apply

Can use undo to undo the last command

**Current relation**

Relation: iris–weka.filters.unsupervised.attribute.Disc...
Instances: 150          Attributes: 5

**Selected attribute**

Name: petallength          Type: Nominal
Missing: 0 (0%)     Distinct: 10     Unique: 0 (0%)

**Attributes**

| No. | Name |
| --- | --- |
| 1 | sepallength |
| 2 | sepalwidth |
| 3 | petallength |
| 4 | petalwidth |
| 5 | class |

| Label | Count |
| --- | --- |
| '(–inf–1.45]' | 23 |
| '(1.45–1.55]' | 14 |
| '(1.55–1.8]' | 11 |
| '(1.8–3.95]' | 13 |
| '(3.95–4.35]' | 14 |
| '(4.35–4.65]' | 15 |
| '(4.65–5.05]' | 18 |

Colour: class (Nom) | Visualize All

Note there are 10 bins
Frequencies not equal but
as close as possible

23  14  11  13  14  15  18  12  17  13

**Status**

OK          Log          x 0

# Weka-2 Review: Preprocessing

- You should now be able to:
  - Import data into Weka
  - Visualize the features and how they relate to class
  - Have an idea of the preprocessing facilities (filters) in Weka and be able to apply them

# Weka-3: Decision Trees

- Weka supports many classification and regression algorithms
- We will go over a DT example is some detail
    - Build a decision tree to classify Iris data set
    - Examine the results
    - Visualize the decision tree
    - Visualize the errors

# Weka-3: Prediction Algorithms

- Weka supports all major classification and regression methods
  - Decision Trees, Rule learners, Nearest Neighbor, Naïve Bayes, Neural Networks, etc.
- Also support ensemble classifiers:
  - Will learn about these later, but combine classifiers
- We first focus on decision trees
- Let's start fresh without the discretization
  - Undo last step or close window and reload "Iris"

Click on Classify Tab

This is just the first classifier that shows up

Again, we have reverted to slides with old version that looks a bit different

Lots of choices, organized under several categories. We will use J48 that can be found under "Trees".

Note that the version of Weka we are using has a somewhat different organization

Click in this area

Accept defaults. Just click OK.

Build and evaluate model using a percentage split instead of 10-fold cross validation. The 66% is for the training set so 34% is for testing.

Click on "More Options" to open this window

Right click on model
Window pops up
then select visualize tree

# Weka-3 Review: Decision Trees

- Covered an example in detail
- You should be able to build a decision tree classifier and understand the results
- You should be able to find the necessary documentation to understand and set the various parameters

# Weka-4: Other Classifiers and Feature Selection

- Briefly explore two more classifiers
  - Artificial Neural Networks
  - Nearest Neighbor
- Demonstrate that with Weka it is easy to run any classifier and can treat like "Black Box"
- Briefly introduce WEKA feature selection

# Artificial Neural Networks (ANNs)

- ANNs covered later in course
- Inspired by the brain
  - Simple processing nodes connected together
  - Each node has a weight and learning occurs by adjusting the weights to get desired output
- ANNs also called Multilayer Perceptrons
- Examples of universal function approximators

# Creating ANN in WEKA for Iris

- While processing Iris data set, choose ANN:
  - In Classify tab click "Choose" button
  - Navigate to "Classifiers" and then "functions"
  - Select "Multilayer Perceptron"
  - Make sure the training/test is still set 66% training
  - Click "Start" to run ANN (use default options)

# Nearest Neighbor Classification

- Nearest Neighbor classifies example based on most similar instance(s)
  - Relators use this method to price your home
  - Also called instance-based learning
  - $k$ used to specify number of nearest neighbors
    - Common values 1, 3, 5
    - Method sometimes referred to as $k$NN
    - Weka uses the IB$k$ algorithm (IB=Instance Based)
  - Also called "Lazy Learning"
    - No work is done up front to build model
    - All work done at classification time

# Creating IBk Model with Weka

- While processing Iris data set, choose IBk:
  - In Classify tab, click "Choose" button
  - Navigate to "Classifiers" and then "Lazy"
  - Select "IBk"
    - Keep all of the default options
  - Make sure the training/test is still set 66% training
  - Click "Start" to run method
    - Results on next slide (uses default k=1)
    - Results for k=5 on slide after that: try it by modifying the option

Note k=5

Note 1 error total for 5NN

# Run a Few Classifiers on Own

- Run Random Forest
  - A tree ensemble method that is under Trees
    - It builds lots of trees using slightly different training examples and features and then combines them
    - We will study them and other ensembles later in course
  - The *numIterations* parameter controls # of trees
- Run Bagging
  - An ensemble (meta learning) method, so find it under classifiers→meta
    - Bagging runs base classifier repeatedly with different training data
    - Run with default classifier and then change classifier to J48

# Weka Feature Selection

- Explore simple feature selection methods
- Use Vote Dataset that has more features
  - Go to Preprocess tab and load using url:
    - https://storm.cis.fordham.edu/~gweiss/data-mining/weka-data/vote.arff
    - Feature selection is needed for methods that cannot handle irrelevant or redundant features
- Experiments:
  - Best first forward search
    - Pick best feature and then add more, but with backtracking to explore more of the space
  - Info gain with ranking
    - Rank the features based on their information gain

## Best First Forward Search for feature subset



1) Click on "Select Attributes"
2) Hit start using the defaults

Results window may only have room for experiment info so scroll down to see actual results (displayed next slide)

# Best First Search Subset Selection

```
=== Attribute Selection on all input data ===

Search Method:
        Best first.
        Start set: no attributes
        Search direction: forward
        Stale search after 5 node expansions
        Total number of subsets evaluated: 85
        Merit of best subset found:    0.729

Attribute Subset Evaluator (supervised, Class (nominal): 17 Class):
        CFS Subset Evaluator
        Including locally predictive attributes

Selected attributes: 3,4,10,11 : 4
                     adoption-of-the-budget-resolution
                     physician-fee-freeze
                     immigration
                     synfuels-corporation-cutback
```

4 attributes
selected of the 17

# Rank Features by Info Gain



1) Modify both defaults as shown
2) Hit Start
3) Scroll to end of results window

# Review Weka-4: Other Classifiers and Feature Selection

- Built and evaluated neural network and nearest neighbor models
  - You should be able to use Weka to run any classification model
- You should understand how to run feature selection using Weka
  - It is okay if you do not understand the details of each method
    - Can look them up in the documentation

# Weka-5: KnowledgeFlow Interface

- Provides a drag and drop interface
  - Each step is a node that you drop onto workspace
  - You connect nodes to create a process flow

- Advantage: supports multiple process flows
  - Can have a path with steps to preprocess data etc.
  - Then branch to build different classifiers

- Other DM tools provide similar interfaces
  - My early DM tools in industry had this interface
    - SAS Enterprise Miner
    - Clementine (now IBM SPSS Modeler)

# Video Demonstrating KnowledgeFlow

- You are not required to try KnowledgeFlow
- Watch first 4 minutes of 7-minute video
    - It does not show how to use mouse to connect nodes, which can be a bit tricky
    - It conveys the basic ideas and power of interface
    - Link to Ian Witten's Youtube video:
        - https://youtu.be/sHSgoVX9z-8
        - Also listed under my Youtube playlist:
            - https://www.youtube.com/playlist?list=PLg2TXNAfR8_ukJY1fEPDMuXV43ZrBBpUh
            - Select "More Data Mining with Weka (1.4 .."

# Review of DM with Weka Module

- You should have installed Weka
- Understand the different interfaces
- Using Explorer interface, be able to:
  - Load any data set
  - Visualize the features
  - Build any classifier and know how to learn about parameters and change them
  - Examine and interpret the results
  - Know how to perform feature selection
- Understand basics of the knowledgeFlow