# CS286 Lab 2: Exploratory data analysis in WEKA

### K. A. Tarnowska

### 31-Aug-2020

## Contents

# 1   Introduction

In this Lab, you apply the concepts of Rough Set Theory to the dataset chosen in Lab 1. In particular, you will describe the dataset as a *decision table*. You will analyze the *decision-making process* that the dataset simulates and *knowledge* represented by the decision table. You will learn how to install and import the dataset into WEKA and perform You will get hands-on experience with Exploratory Data Analysis in WEKA.

**Pre-requisites**   This lab assumes the completion of Lab 1.

# 2   Rough Set Theory [1]

## 2.1   Decision Table

Based on the dataset that you chose for Lab 1, **answer the following questions**:

- How many *objects* does the *decision table* describe?

- What does the *object* represent? (i.e. patient, log entry)

- How many *attributes* describe the objects?

- What are the *conditional* and what attribute(s) is the *decision*?

## 2.2   Knowledge Representation

Based on the dataset that you chose for Lab 1, **answer the following questions**:

- What *knowledge* is represented by the chosen decision table?

- Who is the *expert* involved? (i.e. physician)

- What *decision-making process* does it simulate? (i.e. diagnosis, treatment)

# 3   Introduction to WEKA [2]-[5]

WEKA is one of several well-known free machine learning choices. Not only is the software widely used, but it is complemented by a textbook [2] an online course [4], and extensive YouTube [5] instructional videos. WEKA requires no programming or math skills and the graphical user interface (GUI) is intuitive, in terms of uploading and navigating the software. It was developed at the University of Waikato in New Zealand. It is free & open source and supports Windows, macOS, and Linux platform. It is Java-based and provides for Java API for embedded machine learning.
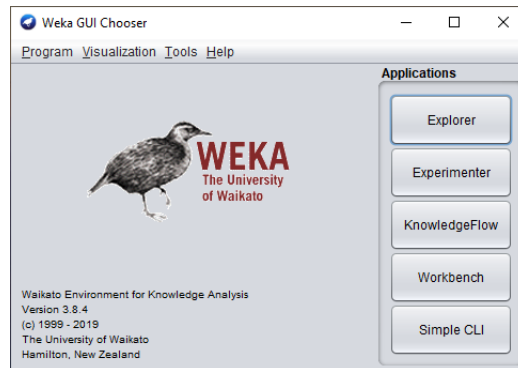
Figure 1: WEKA GUI chooser.

## 3.1 Input

The input can be a csv file, an arff file (attribute-relation file format), which is a text file that includes metadata, xlx, xlsx, json and several others. Data can also be uploaded from a URL or a remote database. WEKA comes with multiple sample datasets, so the user can train on validated data.

## 3.2 User Interface

When the program is opened the choices are (see Figure 1): Explorer (default program), Experimenter (where you can compare model results with t-tests), KnowledgeFlow (Visual workflow), Workbench (combines all of the GUIs into one interface) and SimpleCLI (command-line interface).

**The Explorer** The Explorer option is divided into 5 tabs (see Figure 2): Preprocess (upload, visualize, filter, clean, etc.), Classify (classification and regression algorithms), Cluster (unsupervised), Associate (unsupervised), Select Attributes (ranks importance of variables) and Visualize (creates a master scatter plot so you can see if one variable seems to correlate with another).

## 3.3 Installation

**Download** WEKA from [3]. **Perform installation** steps as described in [3]. Open WEKA Explorer (see Figures 1 and 2).

## 3.4 Loading Data

**Load** the dataset you chose in Lab 1 and described in Section 2. **Include screenshot** showing imported data. See Figure 3 as an example. Explore information displayed in WEKA: relation name, the number of instances, the number of attributes (they all should be visible on the provided screenshot).
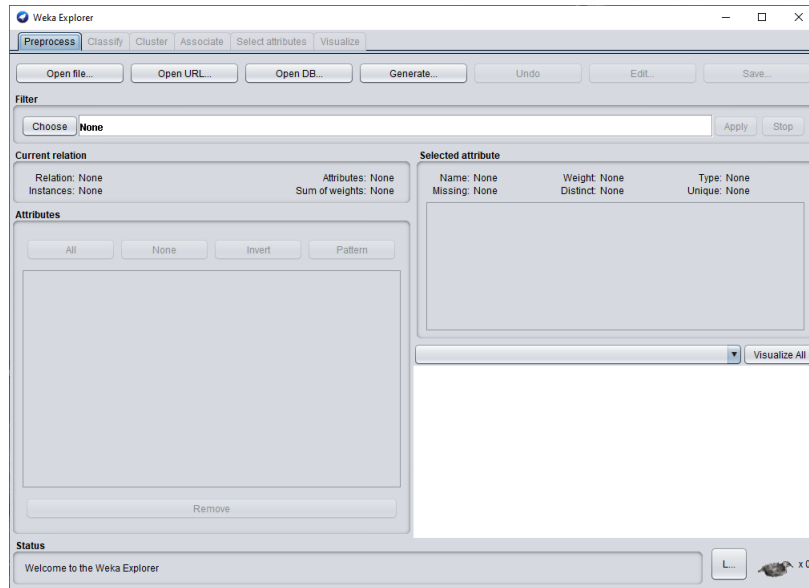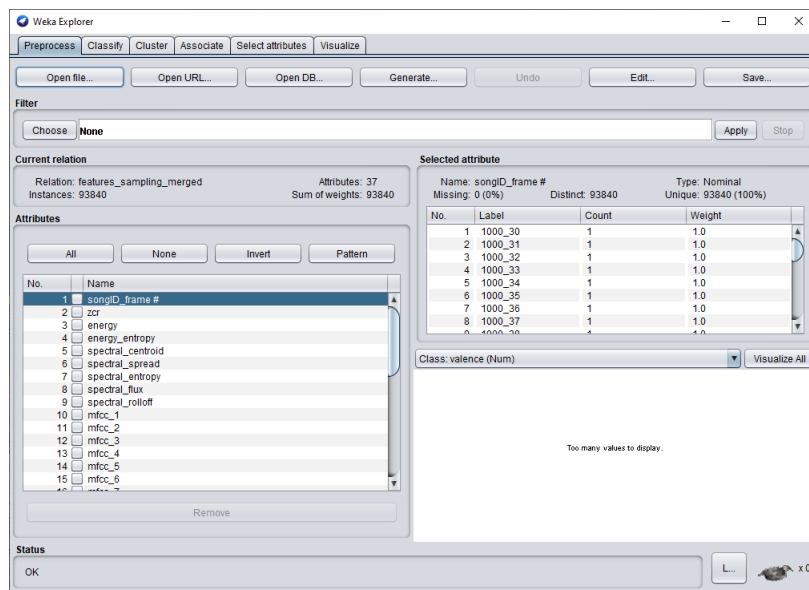
Figure 2: GUI for WEKA Explorer



Figure 3: Example of loading data into WEKA (here, audio dataset)

# 4 Exploratory Data Analysis

**Data analysis** refers to hands-on data exploration and evaluation. **Exploratory analysis** is an approach to analyzing datasets to find previously unknown relationships. Often such analysis involves using various data visualization approaches.

**Statistical vs. Exploratory Data Analysis**  Statistical approaches have "a priori" (presupposed by experience) hypothesis to test. For example, apriori question: Has increasing fee-structure led to decreasing market share" leads to formulating the hypothesis: "Market share has decreased". Next, the analyst applies statistical methods to test the hypothesis. However, not always we have a priori notions about data. In this case, use **Exploratory Data Analysis (EDA)**. It is approach useful for:

- Delving into data

- Examining important interrelationships between attributes

- Identifying interesting subsets of the objects

- Develop an initial idea of possible associations amongst the attributes, as well as between the predictors and the decision.

General goals of EDA are:

- Investigate variables

- Examine *distributions* of *categorical* variables

- Look at *histograms* of *numerical* variables

- Explore relationships among sets of variables (i.e. using *correlation*)

**Variables**  Variables (or attributes) describe the data we are working with. An attribute (or a variable) is a label we give to our data. The data can be *numerical* or *categorical*. **Categorical** variables have values that can be separated into distinct categories (i.e. animal taxonomy: mammals, reptiles, etc.). If we use numbers to represent categories, they become **nominal** variables. However, these numbers are only used to represent the categories, but cannot be used for meaningful mathematical or statistical operations. A variable that is thought not to be controlled or affected by other variables is called an **independent** variable. A variable that depends on other variables is called a **dependent** variable.

**Frequency Distribution**   It is useful to plot a graph showing how many times each value occurs. The common way to visualize frequency distribution is to use a **histogram**. Histograms plot values of observations on the horizontal axis, with a bar showing how many times each value occurred in the dataset. The *normal distribution* is characterized by a bell-shaped curve, where data is distributed symmetrically around the center of all values. There are two ways in which a distribution can deviate from normal:

- Lack of symmetry (called **skew**) - a skewed distribution can be either positively skewed ("to the left") or negatively skewed ("to the right")

- Pointiness (called **kurtosis**) - refers to how "pointy" a distribution is.

**Measures of Centrality and Distribution**   Often one number can tell us enough about a distribution. This is typically a number that points to the "center" of a distribution. There are three measures commonly used: **mean**, **median**, and **mode**. *Mean* is commonly known as an **average**. The *median* is the middle value for a dataset. The *mode* is the most frequently occurring value in a dataset. Distributions come in all shapes and sizes. Simply looking at the central point (mean, median, or mode) may not help in understanding the actual shape of a distribution. Therefore, we often look at the *spread*, or the *dispersion*, of a distribution. The common measures are **range**, **variance**, and **standard deviation**. The **range** is calculated by subtracting the smallest value from the largest value. Sometimes we discard extreme values, and calculate range only for the middle 50% of values. This is known as the **interquartile range**. The *variance* and *standard deviation* indicate how spread out the data points are. To measure the variance, the common method is to pick a center of a distribution, typically the mean, and measure how far each data points from the center. The variance gives us the measure of spread in units squared. The standard deviation is the squared root of the variance, which ensures the measure of average spread is in the same units as the original measure.

**Correlation**   Correlation is a statistical analysis that is used to measure and describe the *strength* and *direction* of the relationship between two variables. Correlation is a simple statistical measure that examines how two variables change together over time. This indicates how closely two variables are related and ranges from -1 (negatively related) to +1 (positively related.) A correlation of 0 indicates no relation between two variables.

## 4.1   Get to Know your Data Set

Graphs, plots, and tables often uncover important relationships in data.

### 4.1.1   Viewing data

**Provide** a screenshot of the subset of objects from your data set presented in the tabular format. You can use the WEKA Explorer *Edit* button to open your

Figure 4: Displaying data instances (objects) in the Viewer of WEKA Explorer.

data in the *Viewer* (see Figure 4 for an example).

### 4.1.2 High-level data analysis

**Provide high-level insights** from your EDA. Describe attributes, their meanings, their types (i.e numerical, categorical, nominal), the decision attribute, the type of decision attribute (i.e. binary, categorical), data missing.

### 4.1.3 Fine-grained attribute analysis

Select an attribute and explore the information provided on the *Selected attribute* section, such as minimum, maximum, mean, standard deviation (if numeric) or frequency analysis (if catgeorical). **Provide a screenshot** with the information displayed on the selected attribute. See Figure 5 as an example of exploring numerical attribute and 6 as an example of exploring nominal (categorical) attribute.

### 4.1.4 Summarization and Visualization

**Provide a graph** (histogram or distribution chart) or at least one chosen attribute (see Figure 7 for an example). **Provide a comment** about a distribution. Use the *Visualize* tab and analyze scatter plots. **Provide comment** if you have discovered any correlations between variables, if so **provide the scatter plot** that supports this.
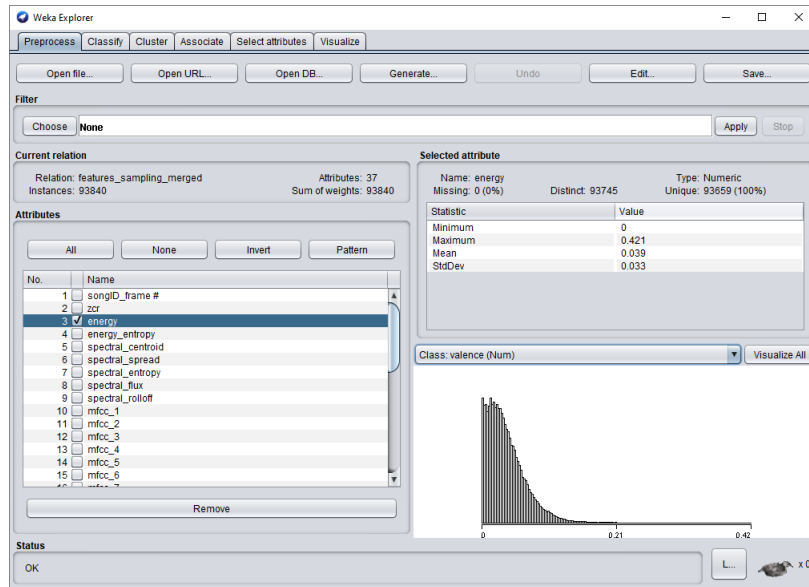
Figure 5: Example of exploring a selected attribute (here, numerical attribute for the energy of the sound), showing positively skewed distribution
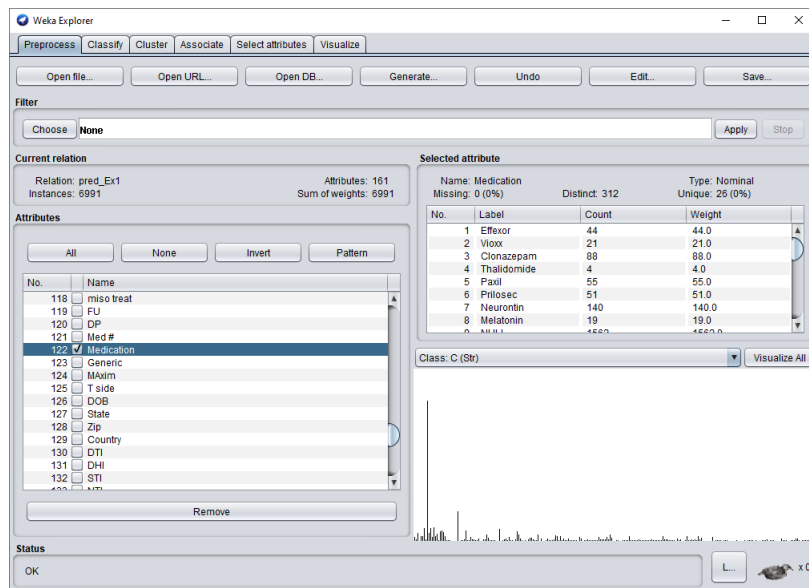


Figure 6: Example of exploring a selected attribute (here, a categorical attribute for the medication of the patient)
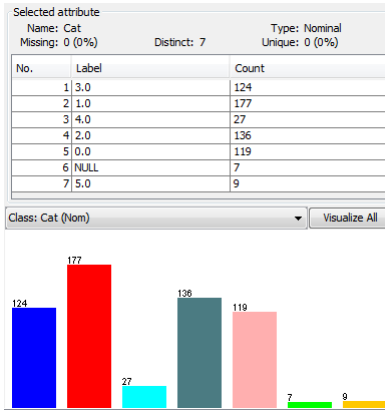
Figure 7: Sample distribution graph of a categorical decision attribute for tinnitus category.

# 5 Deliverables

Lab report uploaded on Canvas (pdf preferred) with embedded screenshots.

# 6 Grading

The lab will be graded according to the rubric in Table 1.

Table 1: Grading rubric

| Rubric | Criteria | Points |
|---|---|---|
| Decision Table | Answered all questions in 2.1 correctly and thoroughly | 2 |
| Knolwedge Representation | Answered all questions in 2.2 correctly and thoroughly | 1 |
| Data loading | Provided proof of successful data loading as described in 3.4 | 1 |
| Data exploration | Provided screenshot of data view as described in 4.1.1 | 1 |
| Data analysis | Provided EDA analysis as described in 4.1.2 | 2 |
| Attribute analysis | Provided screenshot and analysis of the selected attribute as described in 4.1.3 | 1 |
| Visualization | Provided a graph with data distribution and commented as described in 4.1.4 | 1 |
| Correlation analysis | Commented on correlation analysis as described in 4.1.4 | 1 |
| **Total** | | **10** |

# 7 Self-check

Answer the following questions:

1. What is WEKA and what are its main advantages?

2. Describe the data mining process as visualized by tabs on WEKA Explorer.

3. What is EDA and what is it used for?

4. Describe how exploratory analysis differs from statistical analysis.

5. What is a variable? What are the types of variables?

6. Given one example of a numerical and one example of a categorical variable.

7. What are the common techniques used in the exploratory analysis?

8. What are the measures of centrality?

9. What are the measures of distribution?

10. How can we measure the relationship between variables?

# References

[1] Tarnowska KA, Ras ZW, Jastreboff PJ. Knowledge Discovery Approach for Recommendation, Chapter 4 in: Decision Support System for Diagnosis and Treatment of Hearing Disorders. The Case of Tinnitus. Studies in Computational Intelligence. Springer, 2017

[2] Witten IH, Frank E and Hall MA. Introduction to WEKA/The Explorer, Chapter 10 & 11.1, in: Data Mining: Practical Machine Learning Tools and Techniques. Fourth Edition. Morgan Kaufmann. 2017

[3] WEKA 3: Data Mining Software in Java. https://www.cs.waikato.ac.nz/ml/weka/ Accessed April 2, 2019

[4] Future Learn. Data Mining with WEKA. https://www.futurelearn.com/courses/data-mining-withweka Accessed April 3, 2019

[5] YouTube WEKA videos. https://www.youtube.com/results?search_query=WEKA Accessed April 3, 2019