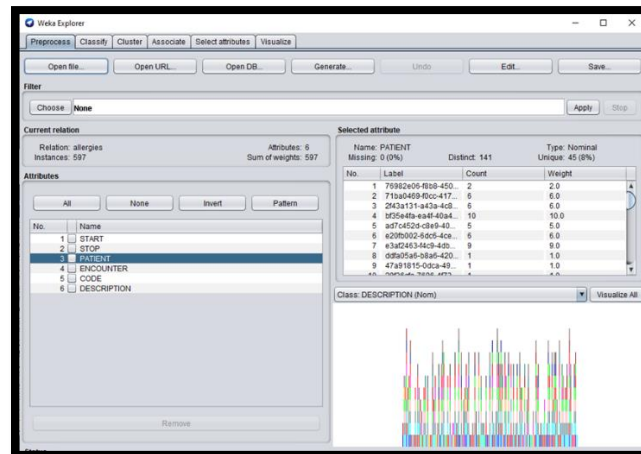




WEKA

Machine learning software to solve data mining problems



Links

- ✓ <https://sourceforge.net/projects/weka/>
- ✓ <https://www.cs.waikato.ac.nz/ml/weka/>
- https://lib.ugent.be/fulltxt/RUG01/000/842/101/RUG01-000842101_2010_0001_AC.pdf

License GNU General Public license ([GPL 3.0 for Weka > 3.7.5](#))
Version 3.8.4
Last Update 12/20/2019
OS Linux, macOS, Windows

Description WEKA (Waikato Environment for Knowledge Analysis) is a collection of machine learning algorithms for solving real-world data mining problems. It is written in Java and runs on almost any platform. The algorithms can either be applied directly to a dataset or called from your own Java code.

It incorporates representation and prescient investigation and displaying strategies, grouping, affiliation, relapse and order.

WEKA is an open source machine learning software that can be accessed through a graphical user interface, standard terminal applications, or a Java API. It is generally used for teaching, research, and industrial applications, contains a built-in tool for standard machine learning tasks, and additionally gives transparent access to well-known toolboxes such as scikit-learn, R, and Deeplearning4j.

WEKA contains a collection of visualization tools and algorithms for data analysis and predictive modeling, together with graphical user interfaces for easy access to these functions.

WEKA supports several standard data mining tasks, more specifically, data preprocessing, clustering, classification, regression, visualization, and feature selection. WEKA's techniques are predicated on the assumption that the data is available as one flat file or relation, where each data point is described by a fixed number of attributes (normally, numeric or nominal attributes, but some other attribute types are also supported). WEKA provides access to SQL databases using Java Database Connectivity and can process the result returned by a database query.

Features

- machine learning
- data mining
- preprocessing
- classification
- regression
- clustering
- association rules
- attribute selection
- experiments
- visualization
- workflow

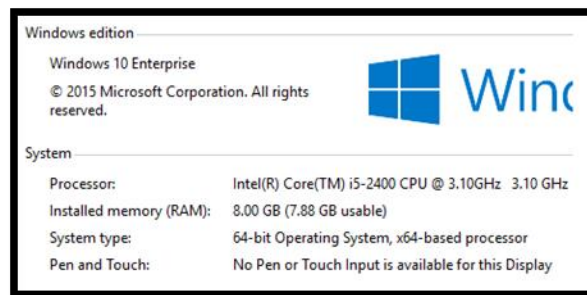
Connectivity / Supported Data Sources & Formats

- Arff, JSON, CSV, xrrf, dat, data, names, and more
- Database using ODBC

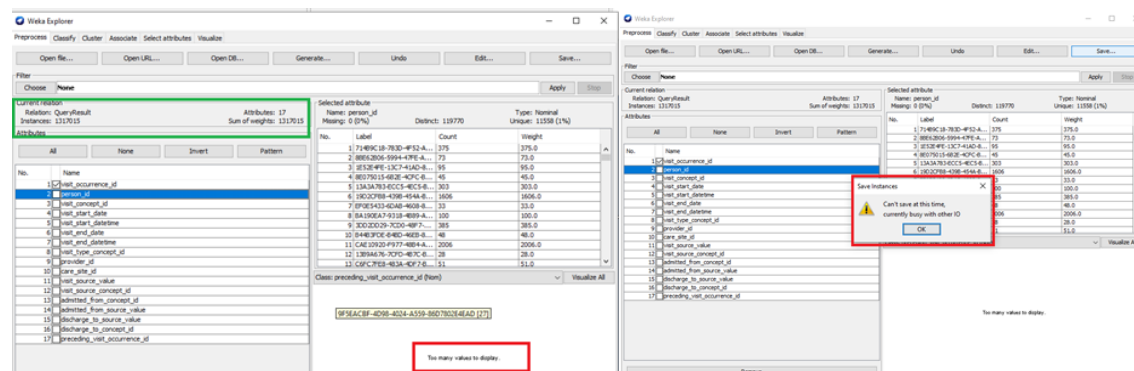
Limitations

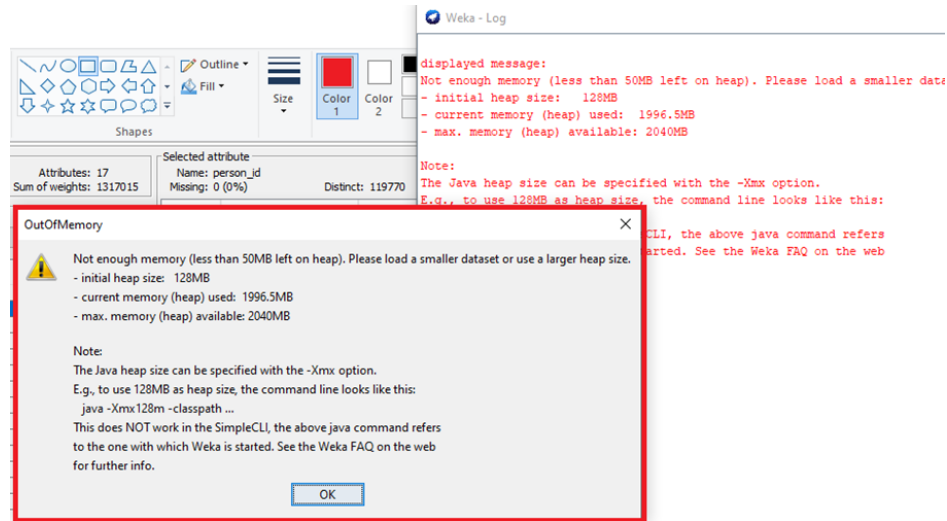
WEKA can only handle small data sets and is not capable of multi-relational data mining, but there is separate software for converting a collection of linked database tables into a single table suitable for processing using WEKA. Another important area that is currently not covered by the algorithms included in the WEKA distribution is sequence modeling.

Performance



Weka loads the 1.3 Million records but does not process the analysis and terminates with "Out-of-Memory" exception as shown in Pic. A machine with more memory and CPU capacity should be used





WEKA DATA ANALYSIS

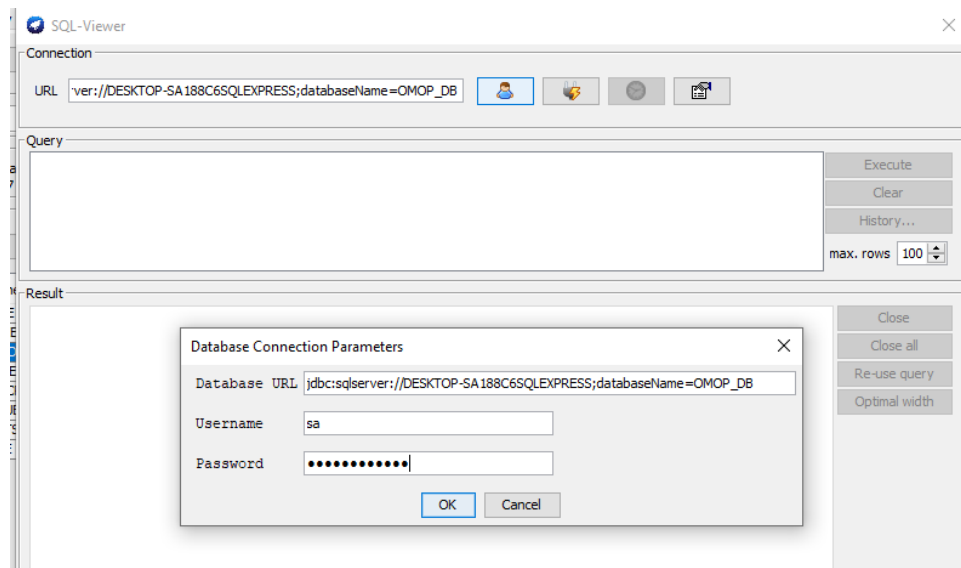
1. Download Weka from "https://waikato.github.io/weka-wiki/downloading_weka/"
2. Install Weka and click "Explorer", Select the "Open File" or "Open DB" options.
3. For CSV File select the file type as "CSV" and browse for the required CSV file.
4. For SQL DB give the connection string in the below format,

jdbc:sqlserver://<ServerName>/<SQLInstanceName>;databaseName=<DBName>

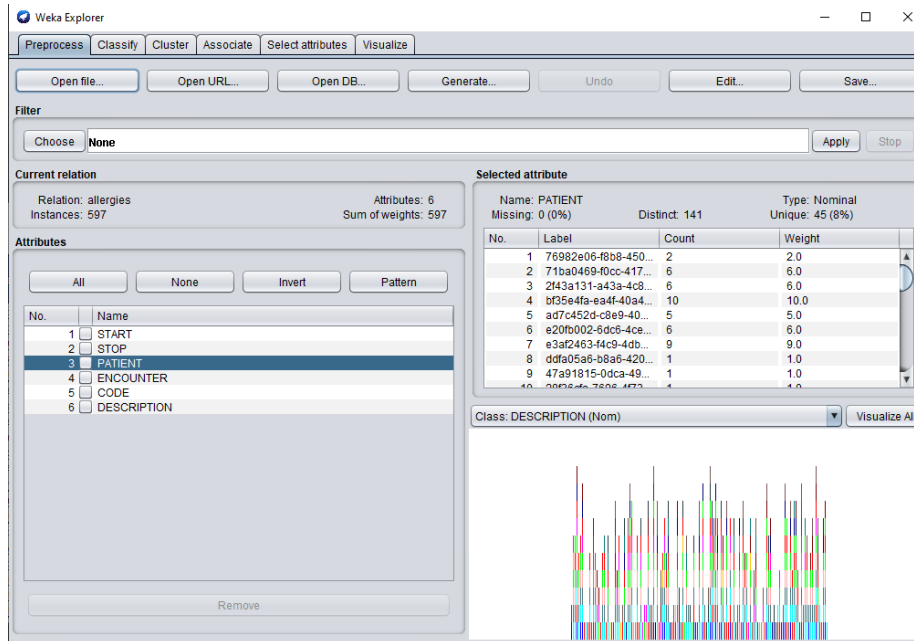
E.g:

jdbc:sqlserver://DESKTOP-SA188C6SQLEXPRESS;databaseName=OMOP_DB

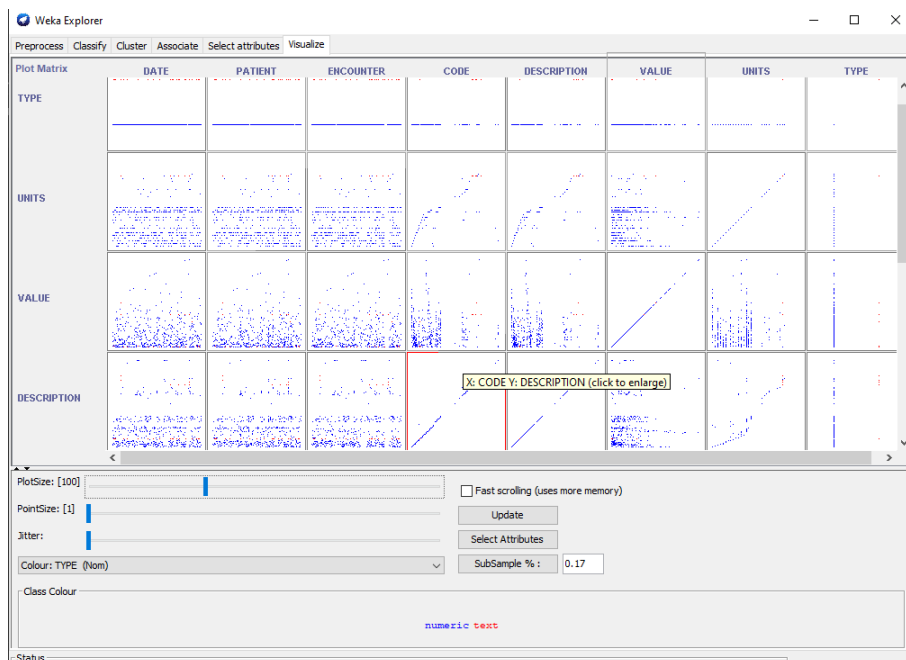
5. Give the required information like username and password.
6. Write a query to execute in the "Query" box and press "Execute".
7. View the data and the statistics as shown below,



SAMPLE WEKA STATISTICS:

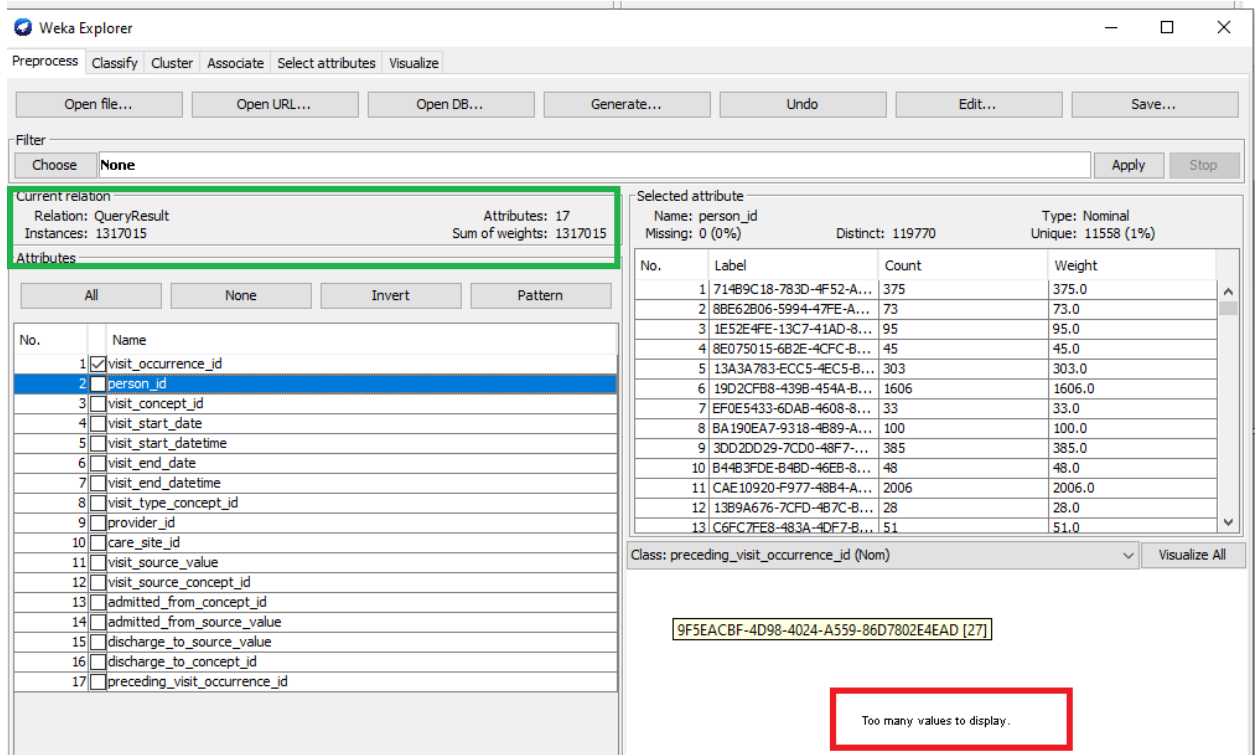


MIN AND MAX VALUES VISUALIZATION:



LIMITATIONS:

- Weka can load the 1,3M records but while trying to process “Too many values” warning is shown and expected results are not shown as shown in below,



Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter: Choose None Apply Stop

Current relation
Relation: QueryResult
Instances: 1317015
Attributes: 17
Sum of weights: 1317015

Attributes

All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> visit_occurrence_id
2	<input checked="" type="checkbox"/> person_id
3	<input type="checkbox"/> visit_concept_id
4	<input type="checkbox"/> visit_start_date
5	<input type="checkbox"/> visit_start_datetime
6	<input type="checkbox"/> visit_end_date
7	<input type="checkbox"/> visit_end_datetime
8	<input type="checkbox"/> visit_type_concept_id
9	<input type="checkbox"/> provider_id
10	<input type="checkbox"/> care_site_id
11	<input type="checkbox"/> visit_source_value
12	<input type="checkbox"/> visit_source_concept_id
13	<input type="checkbox"/> admitted_from_concept_id
14	<input type="checkbox"/> admitted_from_source_value
15	<input type="checkbox"/> discharge_to_source_value
16	<input type="checkbox"/> discharge_to_concept_id
17	<input type="checkbox"/> preceding_visit_occurrence_id

Selected attribute
Name: person_id
Missing: 0 (0%)
Distinct: 119770
Type: Nominal
Unique: 11558 (1%)

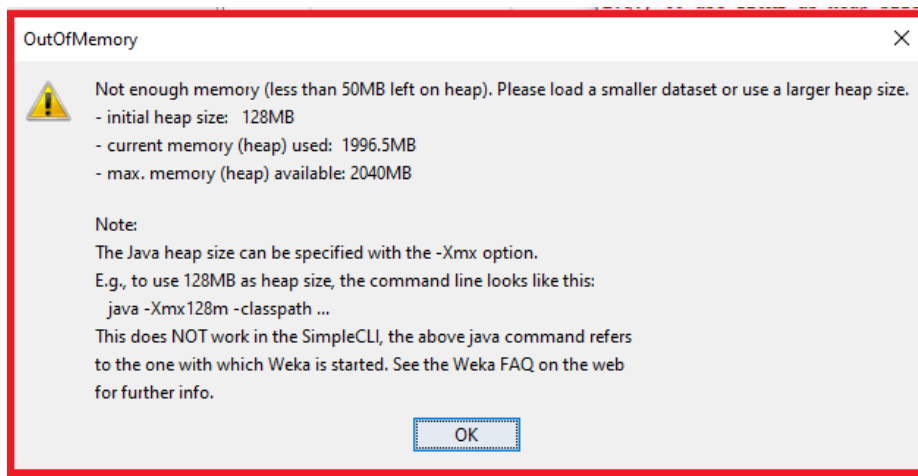
No.	Label	Count	Weight
1	714B9C18-783D-4F52-A...	375	375.0
2	8BE62B06-5994-47FE-A...	73	73.0
3	1E52E4FE-13C7-41AD-8...	95	95.0
4	8E075015-6B2E-4CFC-B...	45	45.0
5	13A3A783-ECC5-4EC5-B...	303	303.0
6	19D2CFB8-439B-454A-B...	1606	1606.0
7	EF0E5433-6DAB-4608-8...	33	33.0
8	BA190EA7-9318-4B89-A...	100	100.0
9	3DD2DD29-7CD0-48F7-...	385	385.0
10	B44B3FDE-B4BD-46EB-8...	48	48.0
11	CAE10920-F977-48B4-A...	2006	2006.0
12	13B9A676-7CFD-4B7C-B...	28	28.0
13	C6FC7FE8-483A-4DF7-B...	51	51.0

Class: preceding_visit_occurrence_id (Nom) Visualize All

9F5EACBF-4D98-4024-A559-86D7802E4EAD [27]

Too many values to display.

- After freezing for some time, the “OutOfMemory” error is displayed.



OutOfMemory

Not enough memory (less than 50MB left on heap). Please load a smaller dataset or use a larger heap size.

- initial heap size: 128MB
- current memory (heap) used: 1996.5MB
- max. memory (heap) available: 2040MB

Note:
The Java heap size can be specified with the -Xmx option.
E.g., to use 128MB as heap size, the command line looks like this:
java -Xmx128m -classpath ...
This does NOT work in the SimpleCLI, the above java command refers to the one with which Weka is started. See the Weka FAQ on the web for further info.

OK