

DATA WAREHOUSING & DATA MINING (01CE0723)

Lab Manual

A.Y. 2025-26

Name : Mahmadaim Kadivar

Er. No. : 92310103056

Semester: 7

Class : TC1

Batch : A

INDEX

Sr. No.	Experiments	Plan Date	Actual Date	Marks	Signature
1.	Explore data mining and data warehousing tools.				
2.	Explore Weka modules: Explorer, Experimenter, KnowledgeFlow, Workbench, Simple CLI. Exploring Explorer module with .csv and .arff files.				
3.	Prepare and analyse “student” dataset, also analyse “student”, “weather.nominal” and “iris” dataset along with editing and visualization.				
4.	Apply Preprocessing techniques on dataset using filters: Remove, ReplaceMissingValues, ReplaceMissingWithUserConstant, ReplaceWithMissingValue, Descritize. Also do the result analysis before and after preprocessing.				
5.	Apply Preprocessing techniques on dataset using filters: NumericToNominal, StringToNominal, NominalToBinary, NumericToNominal. Also do the result analysis before and after preprocessing.				
6.	Demonstration on APRIORI algorithm along with frequent item sets, non-frequent item sets and stron & weak association rules.				
7.	Apply APRIORI algorithm on “weather.nominal” dataset and analyze the results.				
8.	Demonstration on “J48”, “RandomForest” and “NaiveBayes” classification algorithms using test options.				
9.	Apply and analyze “J48”, “RandomForest” and “NaiveBayes” classification algorithms on “weather.nominal” dataset and compare the results.				
10.	Demonstration on prediction algorithms “NaiveBayes” and “Logistic” by creating classification model and “Supplied Test Set” options.				

11.	Apply prediction “NaiveBayes” and “Logistic” by creating classification model and “Supplied Test Set” options on any suitable dataset and compare the results.				
12.	Demonstration on “SimpleKMeans” clustering algorithm using “EuclideanDistance”.				
13.	Apply and analyze “SimpleKMeans” clustering algorithm on suitable dataset, with the observation of “maxIterations” and “numClusters” parameters along with visualization.				
14.	Case study on applications of Data Mining tools and techniques used for Business Intelligence.				

Experiment List

Sr. No.	Title	CO
1.	Explore data mining and data warehousing tools.	CO1, CO2
2.	Explore Weka modules: Explorer, Experimenter, KnowledgeFlow, Workbench, Simple CLI. Exploring Explorer module with .csv and .arff files.	CO2, CO3, CO4, CO5
3.	Prepare and analyse “student” dataset, also analyse “student”, “weather.nominal” and “iris” dataset along with editing and visualization.	CO2
4.	Apply Preprocessing techniques on dataset using filters: Remove, ReplaceMissingValues, ReplaceMissingWithUserConstant, ReplaceWithMissingValue, Descrictize. Also do the result analysis before and after preprocessing.	CO3
5.	Apply Preprocessing techniques on dataset using filters: NumericToNominal, StringToNominal, NominalToBinary, NumericToNominal. Also do the result analysis before and after preprocessing.	CO3
6.	Demonstration on APRIORI algorithm along with frequent item sets, non-frequent item sets and stron & weak association rules.	CO4
7.	Apply APRIORI algorithm on “weather.nominal” dataset and analyze the results.	CO4
8.	Demonstration on “J48”, “RandomForest” and “NaiveBayes” classification algorithms using test options.	CO5
9.	Apply and analyze “J48”, “RandomForest” and “NaiveBayes” classification algorithms on “weather.nominal” dataset and compare the results.	CO5
10.	Demonstration on prediction algorithms “NaiveBayes” and “Logistic” by creating classification model and “Supplied Test Set” options.	CO5
11.	Apply prediction “NaiveBayes” and “Logistic” by creating classification model and “Supplied Test Set” options on any suitable dataset and compare the results.	CO5
12.	Demonstration on “SimpleKMeans” clustering algorithm using “EuclideanDistance”.	CO5
13.	Apply and analyze “SimpleKMeans” clustering algorithm on suitable dataset, with the observation of “maxIterations” and “numClusters” parameters along with visualization.	CO5
14.	Case study on applications of Data Mining tools and techniques used for Business Intelligence.	CO2, CO3, CO4, CO5

Experiment 1

Title: Explore data mining and data warehousing tools.

List of Tools Explored for Data Mining & Data Warehousing:

1. Microsoft SQL Server
2. SAP HANA
3. Oracle Autonomous Data Warehouse + Oracle Data Mining (ODM)

Tool 1: Microsoft SQL Server

- Introduction
Microsoft SQL Server is a relational database management system (RDBMS) developed by Microsoft. It supports both data warehousing and data mining functionalities, especially when used with its integrated services like SQL server database engine, SSIS, SSAS and SSRS.
- Features
Description: - This image shows the Microsoft SQL Server Management Studio (SSMS) interface. The "Object Explorer" panel on the left is expanded and connected to a SQL Server instance named in the image. The explorer lists components such as Databases, Security, Server Objects, Replication, PolyBase, Management, and XEvent Profiler. The central workspace is currently empty, ready for query input or other operations.

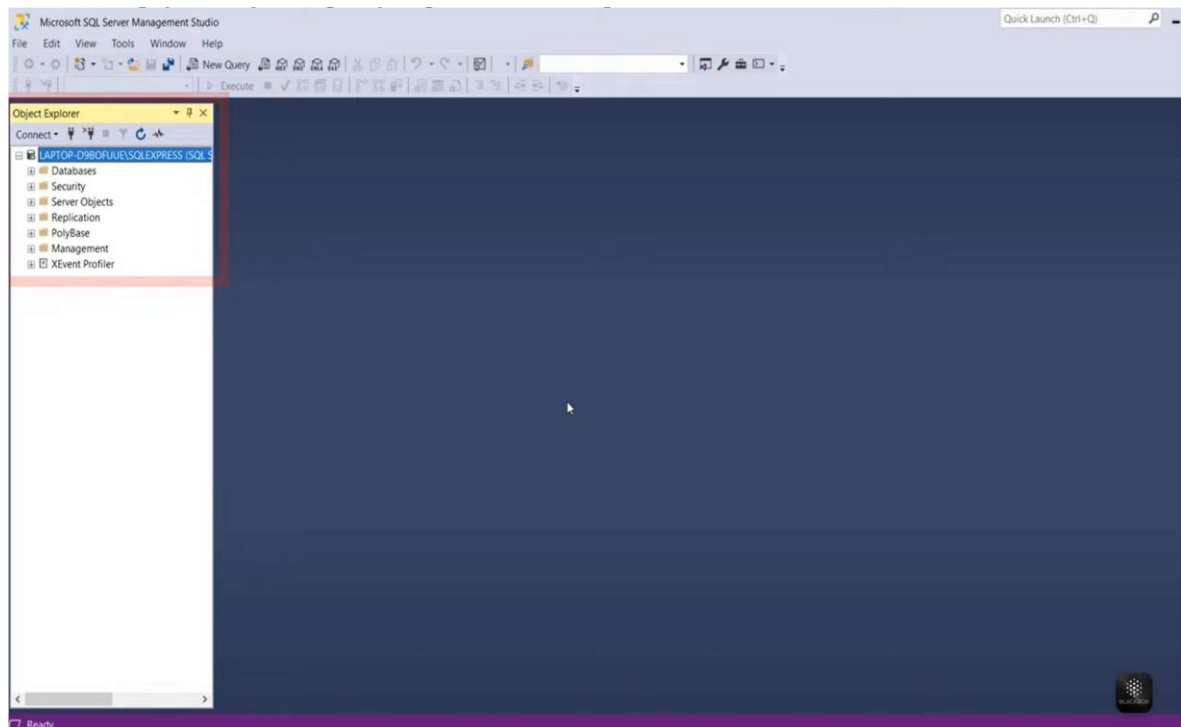


Figure 1.1

Description: - This image shows a simple SQL command has been run to create a table called demo_ssms with just one column named id, which stores integers. Also, the object explorer

displays the connect server and its system databases.

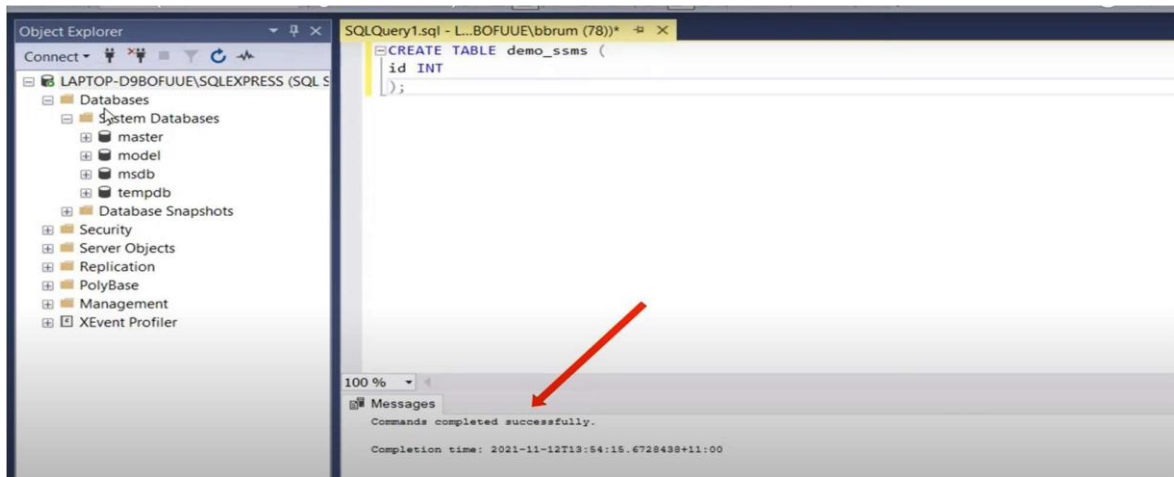


Figure 1.2

- Official website of tool
<https://www.microsoft.com/en-in/sql-server>

Tool 2: SAP HANA

- Introduction
 SAP HANA (High-Performance Analytic Appliance) is an in-memory, column-oriented relational database developed by SAP. It is designed for real-time data processing and analytics, combining transactional and analytical workloads on a single platform.
- Features
 Description: - This image shows the SAP HANA Studio interface. It is divided into three main sections: Systems view, Editor area, and Other views like Error Log and Properties . Each section provides different tools for managing, monitoring, and configuring SAP HANA systems.

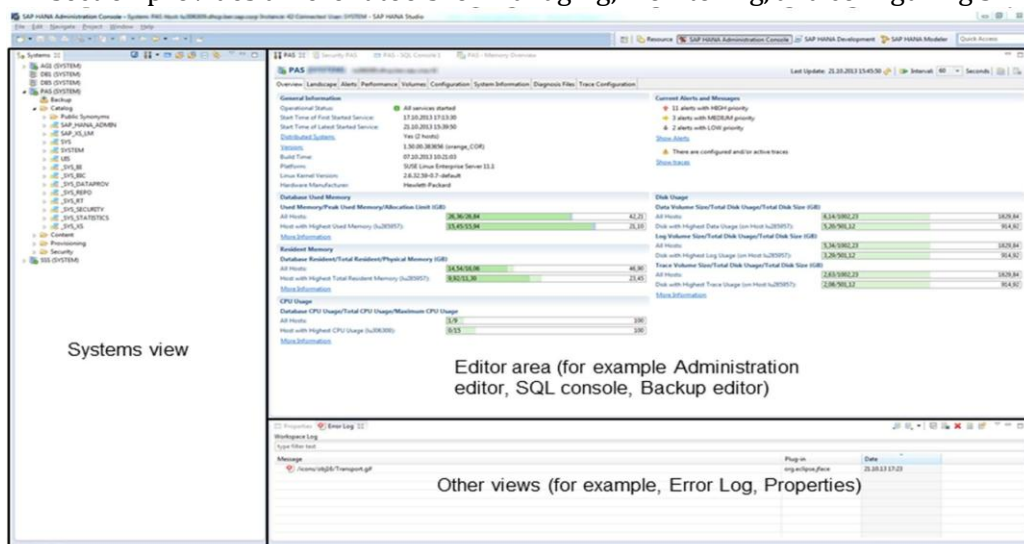


Figure 1.4

- Official website of tool
<https://www.sap.com/products/data-cloud/hana.html>

Tool 3: Oracle Autonomous Data Warehouse + Oracle Data Mining (ODM)

- Introduction
 Oracle Autonomous Data Warehouse (ADW) is a fully managed, cloud-based data warehouse service optimized for analytic workloads, offering high performance, scalability, and automation. Oracle Data Mining (ODM), part of Oracle Advanced Analytics, enables in-database machine learning, allowing users to build predictive models directly within the Oracle Database. Together, they empower businesses to uncover insights and make data-driven decisions efficiently and securely.
- Features
 Description: - This screenshot shows the Oracle Data Miner GUI within Oracle SQL Developer. It features multiple data visualizations such as scatter plots, histograms, and box plots used for exploratory data analysis. The interface supports interactive graphing to help users identify trends and patterns before building predictive models.



Figure 1.5

- Official website of tool
<https://www.oracle.com/in/autonomous-database/>

- **Comparison of all the tools:**

Feature	Oracle Autonomous Data Warehouse + ODM	SAP HANA	Microsoft SQL Server
Vendor	Oracle Corporation	SAP SE	Microsoft
Type	Cloud-based Autonomous Data Warehouse + Advanced Analytics	In-memory database & platform	Relational DBMS
Deployment	Cloud (Autonomous)	On-premise, Cloud, Hybrid	On-premise, Cloud (Azure), Hybrid
Primary Focus	Self-managing data warehouse + integrated data mining	Real-time analytics, in-memory computing	Traditional OLTP/OLAP with BI tools

- **Experiment Outcome:**

The experiment demonstrated that data mining and warehousing tools like Oracle Data Miner (ODM), SAP HANA, and SQL Server offer robust capabilities for data analysis and storage. ODM excels in predictive modeling, SAP HANA provides real-time analytics with in-memory processing, and SQL Server offers integrated services for efficient data warehousing and mining.

Experiment 2

Title: Explore Weka modules: Explorer, Experimenter, KnowledgeFlow, Workbench, Simple CLI.

WEKA History & Introduction:

WEKA (Waikato Environment for Knowledge Analysis) is a popular open-source machine learning software written in Java, developed at the University of Waikato in New Zealand. It provides a comprehensive suite of tools for data pre-processing, classification, regression, clustering, association rules, and visualization. It is widely used in research, education, and industrial applications for data mining and machine learning tasks.

History: -

- In 1993: Development of WEKA began at the University of Waikato, New Zealand, as part of a government-funded research project.
- 1997: The original version of WEKA was a Tcl/Tk-based prototype.
- 1999: WEKA was rewritten entirely in Java, making it platform-independent and more user-friendly.
- 2001: The software became popular after the publication of the book "Data Mining: Practical Machine Learning Tools and Techniques" by Ian H. Witten and Eibe Frank, which used WEKA extensively.
- 2005–present: WEKA has continuously evolved, with contributions from the open-source community and updates from the University of Waikato team

WEKA Applications:

Applications of WEKA are: -

- i. Educational and Research Purposes
- ii. Healthcare and Medical Diagnosis
- iii. Business Intelligence and Marketing
- iv. Text Mining and Natural Language Processing (NLP)
- v. Fraud Detection and Cybersecurity

Modules in WEKA:

1. Explorer
2. Experimenter
3. KnowledgeFlow
4. Workbench
5. Simple CLI

Module 1: Explorer

- Purpose of the Module
The Explorer is one of the main user interfaces in WEKA and serves as a central platform for conducting interactive data analysis and machine learning experiments. It allows

users to easily explore, preprocess, model, and evaluate datasets without needing to write any code.

- Screenshots with description

Description: - This image shows the preprocess tab which is used for loading and preparing data before applying ML algorithm. It's designed to guide you step-by-step through data preprocessing, classification, clustering, and evaluation, making it perfect for both beginners and researchers exploring data.

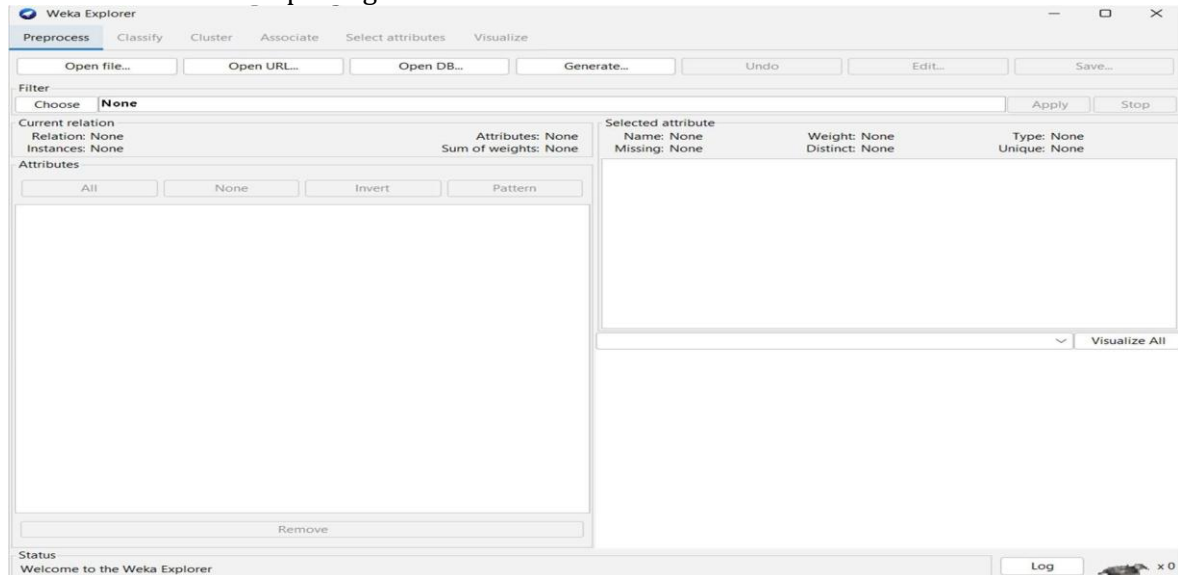


Figure 2.1

- Applications of the module

The WEKA Explorer module is excellent for rapidly evaluating machine learning concepts on actual data. It can be used to clean up your dataset, test various algorithms, such as clustering or decision trees, and quickly assess how well they work.

Module 2: Experimenter

- Purpose of the Module

The WEKA Experimenter makes it easy to test and compare different machine learning models without doing everything by hand. It uses consistent testing methods like cross-validation and even shows you which models perform better with clear, statistical results.

- Screenshots with description

Description: - This image shows the WEKA experimenter environment specially setup tab. It's where you configure your machine learning experiments by choosing datasets, algorithms, and how many times to repeat the tests. The interface lets you easily set up cross-validation, compare models, and organize everything before running the experiments.

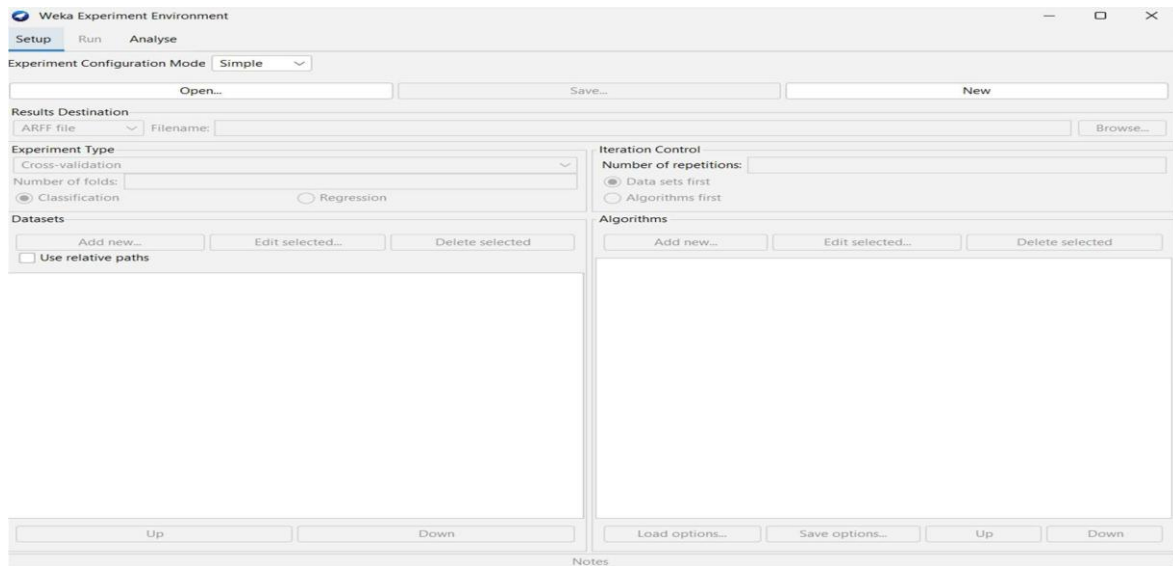


Figure 2.2

- Applications of the module
Its main application is to compare the performance of different algorithms (or their settings) across datasets using consistent evaluation methods like cross-validation. It's especially useful for research, academic projects, or anyone needing to test and validate models in a systematic and repeatable way.

Module 3: KnowledgeFlow

- Purpose of the Module
The KnowledgeFlow module in Weka provides a visual programming environment for designing and executing machine learning workflows. It allows users to connect components like data sources, filters, classifiers, and evaluators in a flowchart-style layout.
- Screenshots with description
Description : - This image shows the WEKA KnowledgeFlow Environment, a visual workspace where you can build machine learning workflows by dragging and connecting components. Instead of writing code, you can design the entire data mining process from loading data and applying filters to training models and visualizing results using a simple flowchart-style interface.

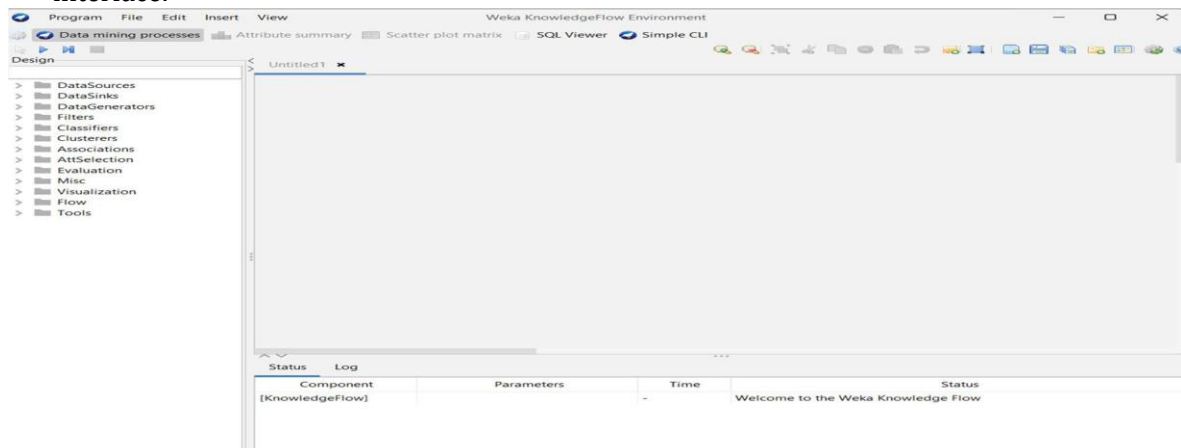


Figure 2.3

- Applications of the module

The WEKA KnowledgeFlow module is mainly used to visually design and automate machine learning workflows. It allows users to connect components like data sources, filters, classifiers, and evaluation tools in a flowchart format.

Module 4: Workbench

- Purpose of the Module

WEKA workbench is the tool used for data mining and ML. It helps to preprocess data, apply various algorithms for classification, regression, clustering and evaluate model performance.

- Screenshots with description

Description: - This image displays the Weka Workbench interface in its initial "Preprocess" state. It's ready for users to load a dataset using options like "Open file..." or "Open URL...". The interface is clean and structured, offering tools for filtering data, viewing attributes, and preparing for further analysis like classification or clustering. No data is loaded yet, as indicated by the empty fields.

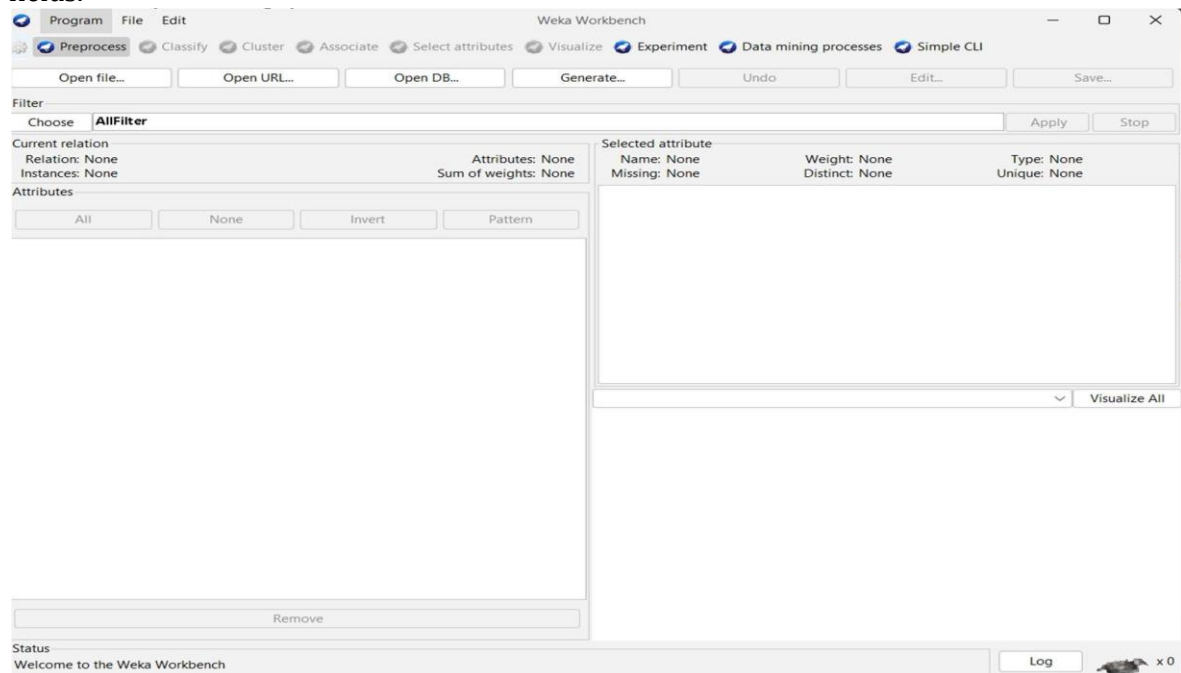


Figure 2.4

- Applications of the module

The Weka Workbench is widely used for data mining, machine learning, and predictive analytics. It helps users preprocess data, build models (e.g., classification, clustering), and evaluate their performance. Common applications include academic research, business intelligence, medical diagnosis, and fraud detection. Its intuitive GUI makes complex analysis accessible without programming.

Module 5: Simple CLI

- **Purpose of the Module**
The Simple CLI in Weka allows users to run machine learning tasks through text commands. It's useful for automation, batch processing, and accessing advanced features. This interface is ideal for experienced users who prefer speed and flexibility over using the graphical interface.
- **Screenshots with descriptions**
Description: - This image shows the Weka SimpleCLI (Command Line Interface), a text-based environment for advanced users to interact with Weka using commands. It allows quick execution of tasks like data preprocessing, model training, and evaluation. The interface supports command history, and shortcuts for ease of use. It's ideal for scripting and automation.



Figure 2.5

- **Applications of the module**
The Simple CLI in Weka is used for automating data mining tasks, running batch processes, and scripting machine learning workflows. It's ideal for large-scale experiments, repeated analyses, and integration into other systems. Researchers and developers often use it to save time and ensure consistency across runs.

Experiment Outcome:

Explorer: Offers an easy-to-use GUI for data preprocessing, classification, clustering, and visualization.

Experimenter: Facilitates systematic comparison of machine learning algorithms through controlled experiments and statistical analysis.

KnowledgeFlow: Provides a visual workflow interface for designing complex machine learning pipelines.

Workbench: Integrates all Weka functionalities into a unified environment for seamless access.

Simple CLI: Allows command-line interaction for quick access to Weka's features, suitable for scripting and automation.

Experiment 3

Title: Prepare and analyse “student” dataset, also analyse “student”, “weather.nominal” and “iris” dataset along with editing and visualization.

• File formats and data types supported by WEKA

- File formats supported by WEKA
 - arff
 - arff.gz
 - bsi
 - csv
 - dat
 - data
 - json
 - json.gz
 - libsvm
 - m
 - names
 - xrff
 - xrff.gz
- Data types supported by WEKA
 - Numeric (Integer and Real), String, Date, and Relational

• Preparation and analysis of “student” dataset

- Dataset Code
@relation student
- @attribute Name string
- @attribute Er_no numeric
- @attribute Email_ID string
- @attribute Mobile_no numeric
- @attribute CGPA numeric
- @attribute Elective_Subject{Cyber,Cloud,Android,Flutter,Network}
- @attribute Class{Tc1,Tc2,Tc3,Tc4,Tc5}
- @attribute Blood_Group{A+,A-,B+,B-,AB+,AB-,O+,O-}
- @attribute Branch{Computer,Electric,Mechanical,Civil}
- @attribute Semester{1,2,3,4,5,6,7,8}
- @attribute Guide_Name string
- @attribute Hobby string
- @attribute Addmission_Year{2020,2021,2022,2023,2024,2025}
- @attribute Gender{Male, Female}
- @attribute Backlogs{0,1,2,3,4,5,6,7,8}
- @attribute Nationality{India,Other}

@attribute DOB date "yyyy-mm-dd"

@attribute Intrested_Job{Professor,Developer,Designer}

@attribute City{Gujrat,Mahrastra,Rajasthan,Tamilnadu}

@attribute State{Morbi,Rajkot,Chennai,Malkapur,Ajmer}

@data

"Mohammadhashim",3055,mohammadhashim45@gmail.com,7096288735,9.5,Cyber,Tc1,O-Computer,7,Shailendra_Chauhan,Coldrinks,2023,Male,0,India,2004-02-08,Developer,Mahrastra,Malkapur

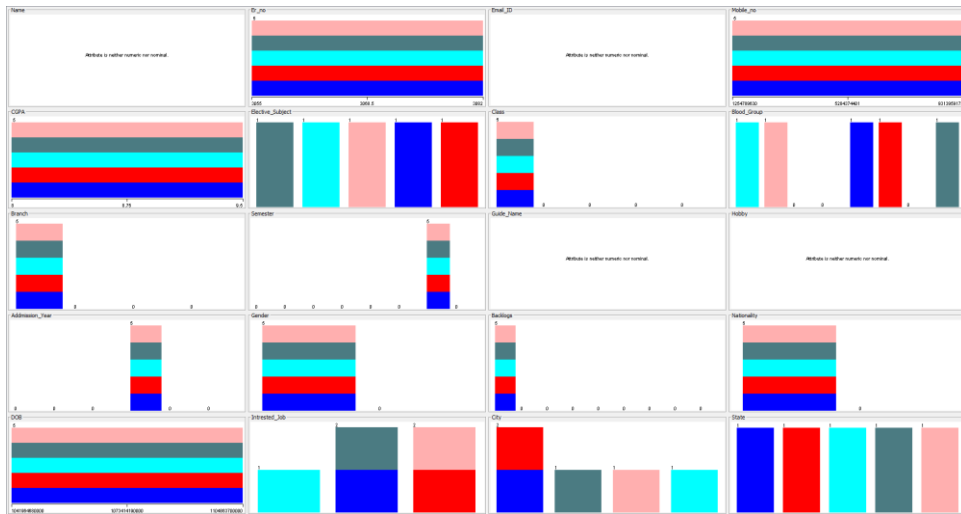
"Mahmadnaim",3056,mahmadnaim07@gmail.com,9313959172,8.5,Cloud,Tc1,A+,Computer,7,Wanglen_Soram,Cricket,2023,Male,0,India,2005-05-05,Professor,Tamilnadu,Chennai

"Samir",3057,Samir18@gmail.com,4785963210,9.0,Android,Tc1,A-Computer,7,Kishan_Makadiya,Sweets,2023,Male,0,India,2003-08-08,Designer,Rajasthan,Ajmer

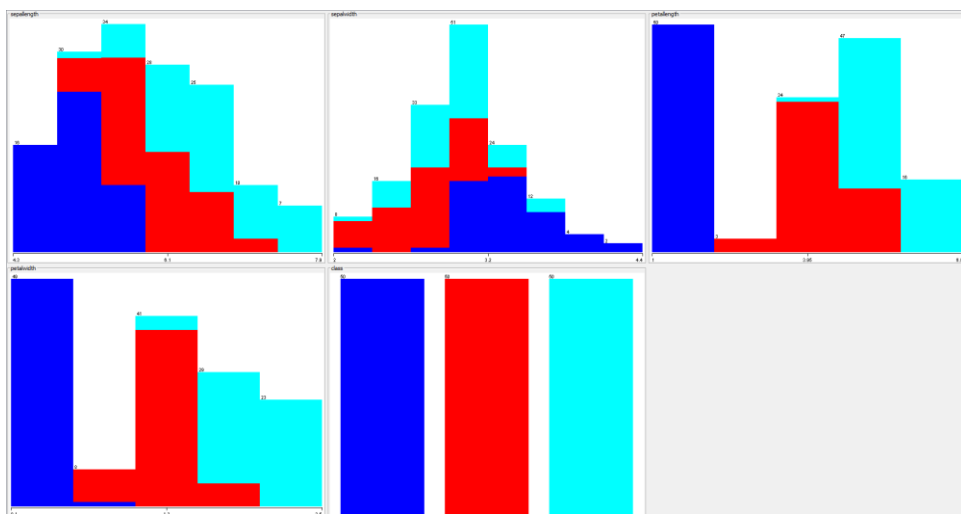
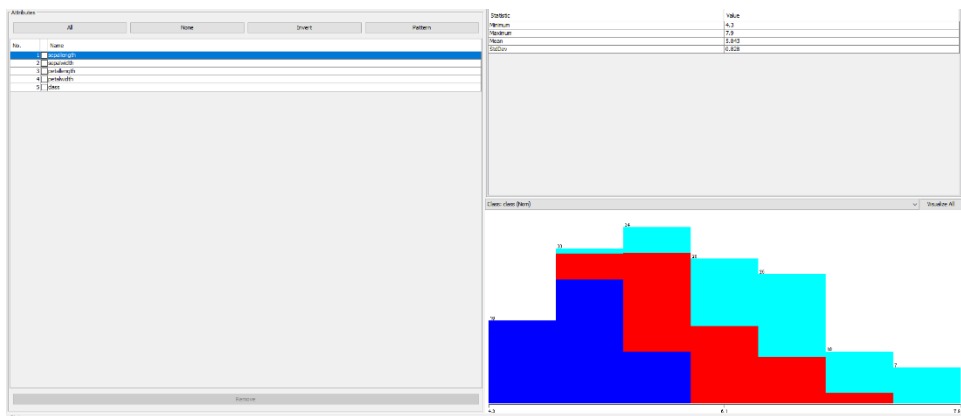
"Yash",3058,yash123@gmail.com,7412589630,8.5,Flutter,Tc1,AB+,Computer,7,Wanglen_Soram,Coldrinks,2023,Male,0,India,2003-09-09,Developer,Gujrat,Morbi

"Pradip",3082,pradip45@gmail.com,1254789630,8.0,Network,Tc1,AB-Computer,7,Kishan_Makadiya,Cricket,2023,Male,0,India,2003-10-10,Designer,Gujrat,Rajkot

Communication				Attributes		Selected attributes		Type: Numeric, Unique: 2 (100%)	
Personnel student				Sum of values: 5		Values: 2 (40%)		Missing: 3 (60%)	
Attributes				Sum of values: 5		Values: 2 (40%)		Missing: 3 (60%)	
No.	Name	Value	Count						
1	Name		1						
2	Age		1						
3	Gender		1						
4	City		1						
5	State		1						
6	Intrested_Job		1						
7	City		1						
8	State		1						
9	Gender		1						
10	Age		1						
11	City		1						
12	State		1						
13	Gender		1						
14	Age		1						
15	City		1						
16	State		1						
17	Gender		1						
18	Age		1						
19	City		1						
20	State		1						
21	Gender		1						
22	Age		1						
23	City		1						
24	State		1						
25	Gender		1						
26	Age		1						
27	City		1						
28	State		1						
29	Gender		1						
30	Age		1						
31	City		1						
32	State		1						
33	Gender		1						
34	Age		1						
35	City		1						
36	State		1						
37	Gender		1						
38	Age		1						
39	City		1						
40	State		1						
41	Gender		1						
42	Age		1						
43	City		1						
44	State		1						
45	Gender		1						
46	Age		1						
47	City		1						
48	State		1						
49	Gender		1						
50	Age		1						
51	City		1						
52	State		1						
53	Gender		1						
54	Age		1						
55	City		1						
56	State		1						
57	Gender		1						
58	Age		1						
59	City		1						
60	State		1						
61	Gender		1						
62	Age		1						
63	City		1						
64	State		1						
65	Gender		1						
66	Age		1						
67	City		1						
68	State		1						
69	Gender		1						
70	Age		1						
71	City		1						
72	State		1						
73	Gender		1						
74	Age		1						
75	City		1						
76	State		1						
77	Gender		1						
78	Age		1						
79	City		1						
80	State		1						
81	Gender		1						
82	Age		1						
83	City		1						
84	State		1						
85	Gender		1						
86	Age		1						
87	City		1						
88	State		1						
89	Gender		1						
90	Age		1						
91	City		1						
92	State		1						
93	Gender		1						
94	Age		1						
95	City		1						
96	State		1						
97	Gender		1						
98	Age		1						
99	City		1						
100	State		1						



- **Analysis of “iris” dataset**



- Analysis of “weather.nominal” dataset

