| Name | Shubham Golwal |
| --- | --- |
| UID no. | 2020300015 |
| Experiment No. | 07 |

| AIM: | Build a Data Warehouse and Explore WEKA Implementation |
| --- | --- |
| **Program 1** | |
| **PROBLEM STATEMENT :** | The objective of this exp is to build a data warehouse using breast cancer datasets and explore WEKA implementation for data analysis. The expt involves data preprocessing, integration, and transformation to build the data warehouse and applying various data mining algorithms using WEKA. The expt is expected to generate insights for breast cancer diagnosis and treatment and contribute to the improvement of healthcare outcomes. |
| **Theory :** | ### What is Weka Tool?<br>Weka is a popular open-source machine learning software tool that provides a wide range of algorithms for data preprocessing, classification, clustering, and regression. It is widely used in both academia and industry for data analysis, data mining, and machine learning tasks.<br><br>In addition to its powerful features, Weka also provides extensive documentation to help users learn how to use the tool effectively. The documentation includes a user manual, a developer guide, and a set of online tutorials that cover various aspects of Weka's functionality.<br>The user manual provides an overview of Weka's features and explains how to install and use the tool. It also includes detailed descriptions of Weka's algorithms, including their underlying mathematical models and their use cases.<br><br>The developer guide provides information on Weka's architecture, including the various modules and components that make up the tool. It also includes information on how to extend Weka's functionality by adding new algorithms or modules.<br><br>The online tutorials cover various topics, including data preprocessing, classification, clustering, and regression. Each tutorial provides step-by-step |

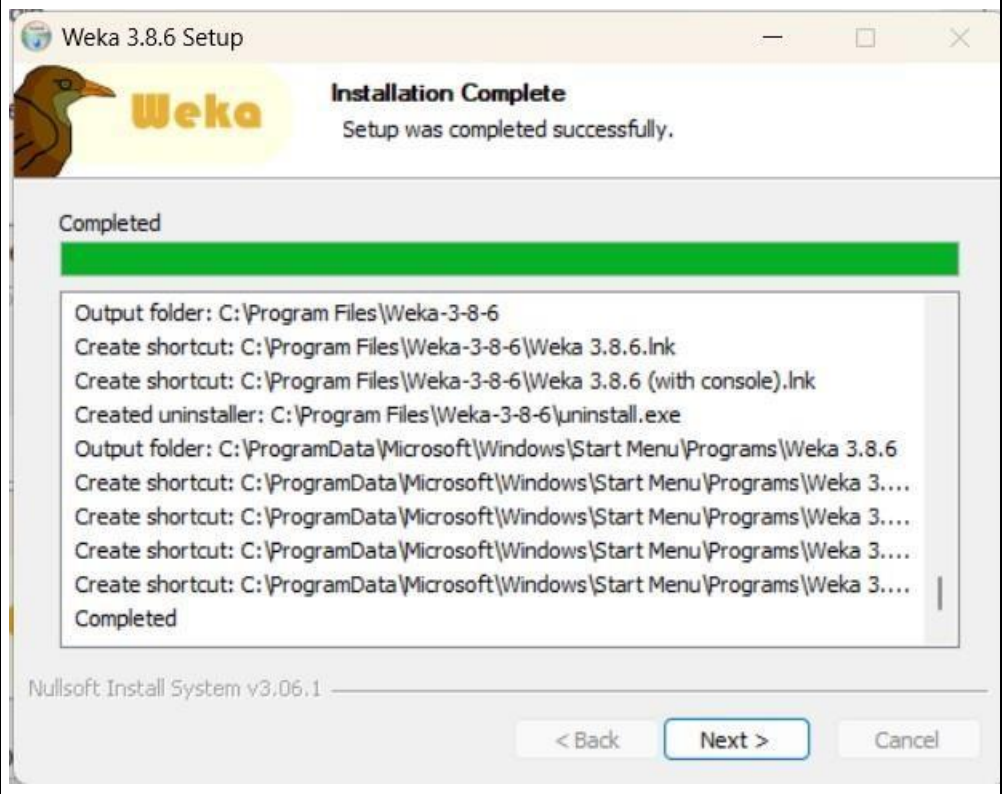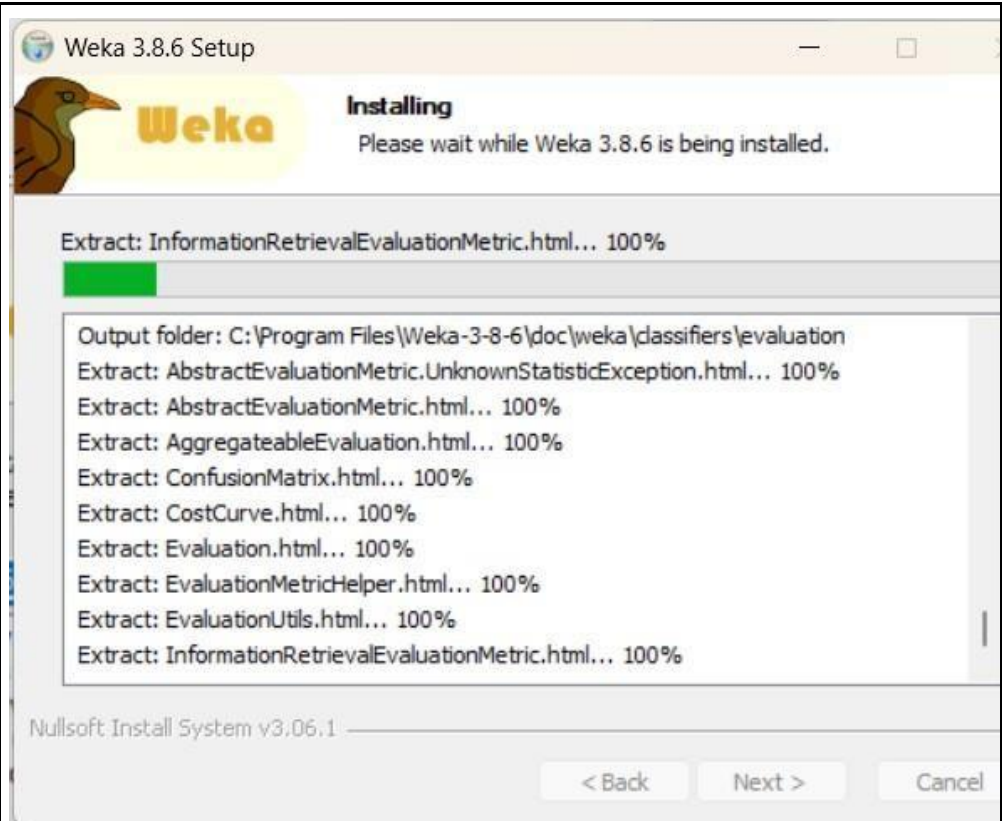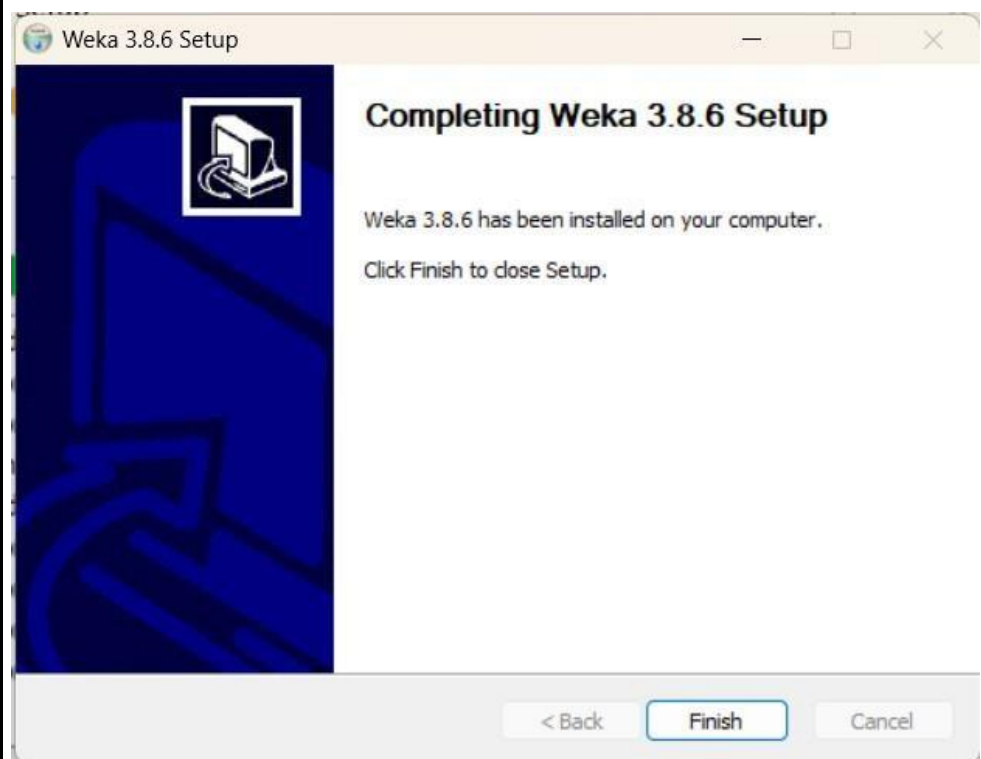| | |
|---|---|
| | instructions on how to use Weka's features to solve a particular problem, and includes sample data and code snippets to help users get started. |
| | Overall, Weka's extensive documentation makes it an accessible and user-friendly tool for data analysis and machine learning tasks. |
| **Installation** |  |

**Weka 3.8.6 Setup** — License Agreement

Please review the license terms before installing Weka 3.8.6.

Press Page Down to see the rest of the agreement.

> GNU GENERAL PUBLIC LICENSE
> Version 3, 29 June 2007
>
> Copyright (C) 2007 Free Software Foundation, Inc. <http://fsf.org/>
> Everyone is permitted to copy and distribute verbatim copies
> of this license document, but changing it is not allowed.
>
> Preamble
>
> The GNU General Public License is a free, copyleft license for
> software and other kinds of works.

If you accept the terms of the agreement, click I Agree to continue. You must accept the agreement to install Weka 3.8.6.

Nullsoft Install System v3.06.1

[< Back] [I Agree] [Cancel]



**Weka 3.8.6 Setup** — Choose Components

Choose which features of Weka 3.8.6 you want to install.

Check the components you want to install and uncheck the components you don't want to install. Click Next to continue.

Select the type of install: **Full**

Or, select the optional components you wish to install:

☑ Associate Files

Description
Position your mouse over a component to see its description.

Space required: 345.2 MB

Nullsoft Install System v3.06.1

[< Back] [Next >] [Cancel]

Weka 3.8.6 Setup

**Choose Install Location**
Choose the folder in which to install Weka 3.8.6.

Setup will install Weka 3.8.6 in the following folder. To install in a different folder, click Browse and select another folder. Click Next to continue.

Destination Folder

C:\Program Files\Weka-3-8-6          Browse...

Space required: 345.2 MB
Space available: 29.6 GB

Nullsoft Install System v3.06.1

< Back    Next >    Cancel



Weka 3.8.6 Setup

**Choose Start Menu Folder**
Choose a Start Menu folder for the Weka 3.8.6 shortcuts.

Select the Start Menu folder in which you would like to create the program's shortcuts. You can also enter a name to create a new folder.

Weka 3.8.6

Accessibility
Accessories
Administrative Tools
Anaconda3 (64-bit)
Android Studio
Dell
Discord Inc
EPSON Projector
Figma, Inc
Git
HP

☐ Do not create shortcuts

Nullsoft Install System v3.06.1

< Back    Install    Cancel

**Weka 3.8.6 Setup**  — □ ⊠

**Installing**
Please wait while Weka 3.8.6 is being installed.

Extract: InformationRetrievalEvaluationMetric.html... 100%

```
Output folder: C:\Program Files\Weka-3-8-6\doc\weka\classifiers\evaluation
Extract: AbstractEvaluationMetric.UnknownStatisticException.html... 100%
Extract: AbstractEvaluationMetric.html... 100%
Extract: AggregateableEvaluation.html... 100%
Extract: ConfusionMatrix.html... 100%
Extract: CostCurve.html... 100%
Extract: Evaluation.html... 100%
Extract: EvaluationMetricHelper.html... 100%
Extract: EvaluationUtils.html... 100%
Extract: InformationRetrievalEvaluationMetric.html... 100%
```

Nullsoft Install System v3.06.1

< Back    Next >    Cancel

---

**Weka 3.8.6 Setup**  — □ ✕

**Installation Complete**
Setup was completed successfully.

Completed

```
Output folder: C:\Program Files\Weka-3-8-6
Create shortcut: C:\Program Files\Weka-3-8-6\Weka 3.8.6.lnk
Create shortcut: C:\Program Files\Weka-3-8-6\Weka 3.8.6 (with console).lnk
Created uninstaller: C:\Program Files\Weka-3-8-6\uninstall.exe
Output folder: C:\ProgramData\Microsoft\Windows\Start Menu\Programs\Weka 3.8.6
Create shortcut: C:\ProgramData\Microsoft\Windows\Start Menu\Programs\Weka 3....
Create shortcut: C:\ProgramData\Microsoft\Windows\Start Menu\Programs\Weka 3....
Create shortcut: C:\ProgramData\Microsoft\Windows\Start Menu\Programs\Weka 3....
Create shortcut: C:\ProgramData\Microsoft\Windows\Start Menu\Programs\Weka 3....
Completed
```

Nullsoft Install System v3.06.1

< Back    Next >    Cancel

| | |
|---|---|
| |  Weka 3.8.6 Setup<br><br>**Completing Weka 3.8.6 Setup**<br><br>Weka 3.8.6 has been installed on your computer.<br><br>Click Finish to close Setup.<br><br>< Back   Finish   Cancel |
| **Output** | **Weka Opened: The landing page of Weka software**<br><br><br><br>**Explorer tab: It allows datasets and the predictions of Classifiers and Clutterers to be visualized in two dimensions.** |

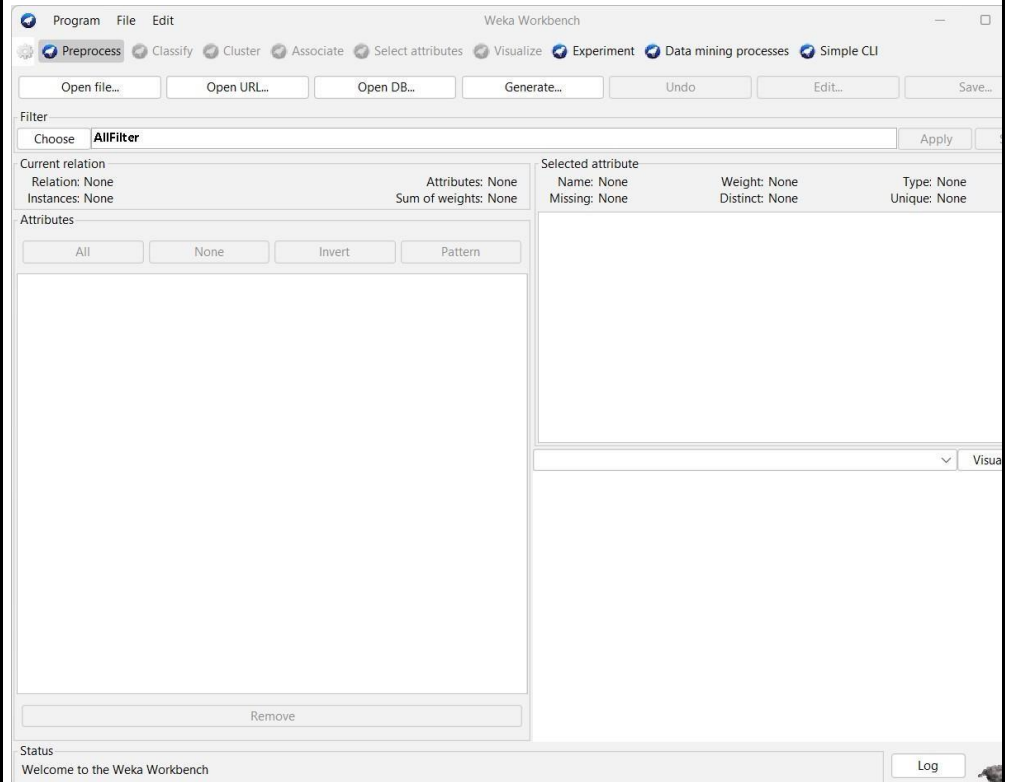**Experimenter tab: The experimenter configures the test options for you with sensible defaults.**

Weka Experiment Environment

Setup    Run    Analyse

Experiment Configuration Mode  Simple

Open...    Save...    New

Results Destination

ARFF file    Filename:    Browse...

Experiment Type                          Iteration Control
Cross-validation                          Number of repetitions:
Number of folds:                          ○ Data sets first
○ Classification    ○ Regression          ○ Algorithms first

Datasets                                  Algorithms
Add new...  Edit selected...  Delete selected    Add new...  Edit selected...  Delete selected
☐ Use relative paths

Up    Down    Load options...  Save options...  Up    Down

Notes

**KnowledgeFlow tab: It presents a "data-flow" inspired interface to Weka**

Program  File  Edit  Insert  View         Weka KnowledgeFlow Environment

○ Data mining processes  Attribute summary  Scatter plot matrix  SQL Viewer  Simple CLI

Design                                    Untitled1

> DataSources
> DataSinks
> DataGenerators
> Filters
> Classifiers
> Clusterers
> Associations
> AttSelection
> Evaluation
> Misc
> Visualization
> Flow
> Tools

Status    Log

Component    Parameters    Time    Status
[KnowledgeFlow]                -        Welcome to the Weka Knowledge Flow

**Workbench tab It contains a collection of data pre-processing tools and machine learning algorithms wrapped in an easy-**

**to-use graphical interface.**



## Simple CLI

**WEKA DATASETS:**



Choosing dataset - We have chosen a breast cancer.arff dataset which contains 3 attributes with yes or no questions and 7 continuous regression matrix.

The breast cancer dataset contains ten attributes that describe various characteristics of breast cancer patients and their tumors. Here's a brief description of each attribute:

1. age: The age of the patient at the time of diagnosis.
2. menopause: The menopausal status of the patient at the time of diagnosis.
3. tumor-size: The size of the tumor, measured in centimeters.
4. inv-nodes: The number of axillary lymph nodes involved in the cancer.
5. node-caps: Whether or not the cancer has spread to the lymph node capsule.
6. deg-malig: The degree of malignancy of the tumor, ranging from 1 (low) to 3 (high).
7. breast: The affected breast (left or right).
8. breast-quad: The quadrant of the breast where the tumor is located.
9. irradiat: Whether or not the patient received radiation therapy.
10. Class: The classification of the tumor as benign or malignant.

These attributes provide important information for assessing the severity and treatment options for breast cancer. Age, menopause, and tumor size can provide insights into the progression of the disease, while inv-nodes, node-caps, and deg-malig can help determine the stage and aggressiveness of the cancer. The breast and breast-quad attributes can help locate the tumor, while irradiat can indicate the type of treatment received. Finally, the Class attribute is the target variable that is used to classify the tumor as either benign or malignant.



## PREPROCESSING DATA:

## Applying apriori algorithm:



```
=== Run information ===

Scheme:      weka.associations.Apriori -N 10 -T 0 -C 0.5 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1
Relation:    breast-cancer
Instances:   286
Attributes:  10
             age
             menopause
             tumor-size
             inv-nodes
             node-caps
             deg-malig
             breast
             breast-quad
             irradiat
             Class
=== Associator model (full training set) ===


Apriori
=======

Minimum support: 0.5 (143 instances)
Minimum metric <confidence>: 0.5
Number of cycles performed: 10

Generated sets of large itemsets:

Size of set of large itemsets L(1): 6

Size of set of large itemsets L(2): 6

Size of set of large itemsets L(3): 4

Size of set of large itemsets L(4): 1

Best rules found:

 1. inv-nodes=0-2 irradiat=no Class=no-recurrence-events 147 ==> node-caps=no 145     <conf:(0.99)> lift:(1.27) lev:(0.11) [30] conv:(10.97)
 2. inv-nodes=0-2 irradiat=no 183 ==> node-caps=no 177     <conf:(0.97)> lift:(1.25) lev:(0.12) [34] conv:(5.85)
```
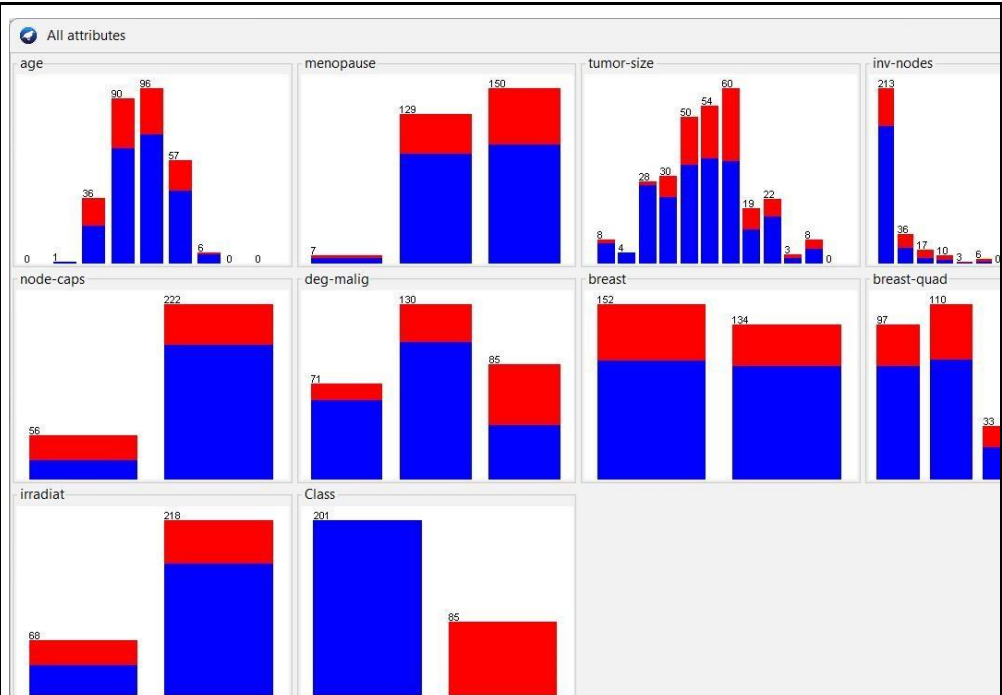
## VISUALISING THE DATSET:

**Conclusion:** Through this experiment, I have gained knowledge and skills in using Weka for data visualization and association. I have become familiar with the Filtered Associator algorithm, which enables efficient exploration of large datasets. I have also explored the command-line interface (CLI) of Weka, which provides more flexibility for advanced users. Overall, I have gained a better understanding of data analysis and exploration using Weka.