# PRACTICAL FILE

## DATA MINING AND PREDICTIVE ANALYSIS LAB

## MLE6806



Submitted By:                                         Submitted To:

Khushi Tanwar                                        Dr. Shivani Sharma

I.B.Tech + M.Tech (AIML)

A501132620001

# **INDEX**

## PROGRAM 1:

To understand the basic features of Data Warehousing.

## THEORY:

A data warehouse is a centralized repository for storing and managing large amounts of data from various sources for analysis and reporting. It is optimized for fast querying and analysis, enabling organizations to make informed decisions by providing a single source of truth for data. Data warehousing typically involves transforming and integrating data from multiple sources into a unified, organized, and consistent format. Below are major **characteristics** of data warehouse :



- **Subject-oriented:** A data warehouse focuses on specific themes like sales or marketing, delivering information tailored to these themes rather than current operations. It eliminates irrelevant data to facilitate precise decision-making.

- **Integrated:** Data integration ensures that data from different sources is combined into a reliable format, allowing effective analysis. Integration involves establishing shared entities and adhering to consistent naming conventions and encoding structures.

- **Time-Variant:** Data in a warehouse is maintained over different time intervals, providing historical perspectives for analysis. Once stored, data cannot be modified, ensuring consistency, and facilitating analysis over time.

- **Non-Volatile:** Data in a warehouse is permanent and not erased or updated when new data is inserted. It is read-only and refreshed at intervals, enabling analysis of historical data without the need for transaction processing or concurrency control.

**Functions of Data warehouse:**

- **Data Consolidation:** Combining multiple data sources into a single repository in a data warehouse for consistency.
- **Data Cleaning:** Identifying and removing errors and inconsistencies from data sources to ensure accuracy.
- **Data Integration:** Combining data from various sources into a unified repository with consistent formatting.
- **Data Storage:** Storing large volumes of historical data for easy access and analysis.
- **Data Transformation:** Cleaning and transforming data to remove inconsistencies and irrelevant information.
- **Data Analysis:** Analysing and visualizing data to extract insights and inform decision-making.
- **Data Reporting**: Providing reports and dashboards for different departments and stakeholders.
- **Data Mining:** Extracting patterns and trends from data to support strategic planning.
- **Performance Optimization:** Optimizing data warehouse systems for fast querying and analysis.

## PROGRAM 2:

Explore WEKA Data Mining/Machine Learning Toolkit:

    (c) Downloading and/or installation of WEKA data mining toolkit.

    (d) Understand the features of WEKA tool kit such as Explorer, Knowledge flow interface, Experimenter, command-line interface.

    (e) Navigate the options available in the WEKA (deselect attributes panel, preprocess panel, classify panel, cluster panel, associate panel and visualize).

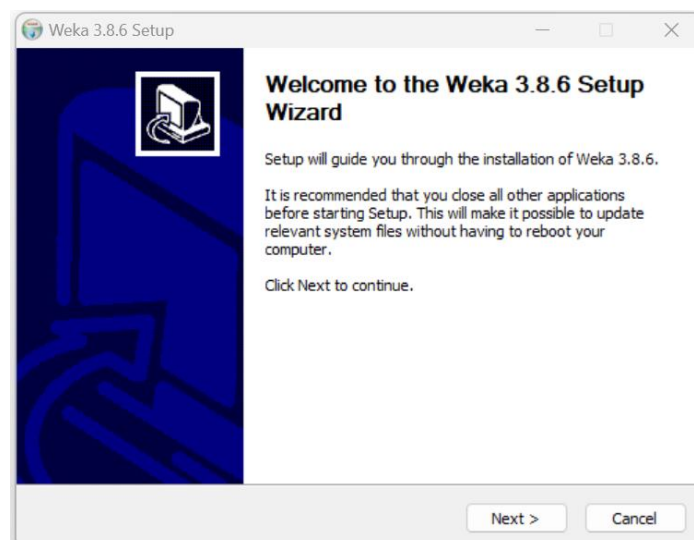    (f) Study the ARFF file format.

## THEORY:

(a) WEKA - an open-source software provides tools for data preprocessing, implementation of several Machine Learning algorithms, and visualization tools so that you can develop machine learning techniques and apply them to real-world data mining problems. How to Download :-
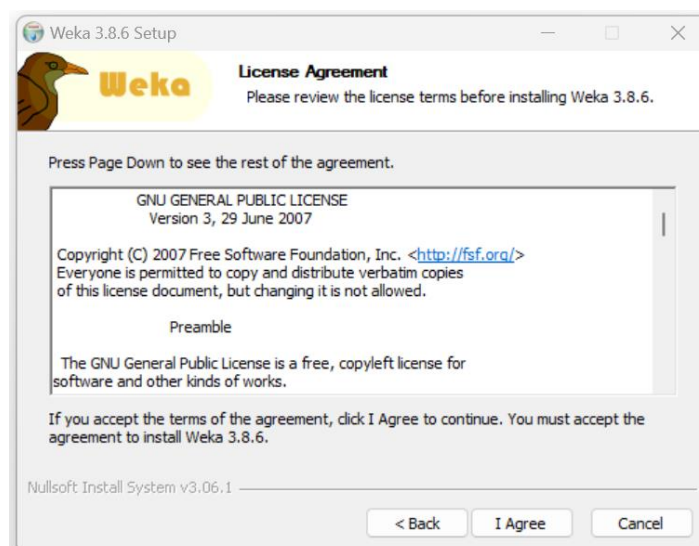
1. **Download the software**(Weka Tool Kit) from the website.

2. The **Java is mandatory** for installation of WEKA so if you have already Java on your machine then download only WEKA else download the software with JVM.

3. Then open the file location and **double click** on the file. Click on **Run**.



4. Click **Next** to continue.

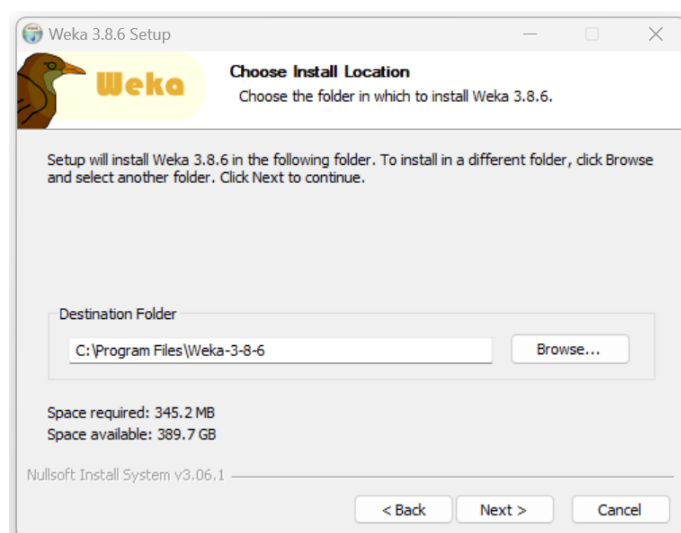5. If you accept the terms of the agreement, click **I Agree** to continue. You must accept the agreement to install Weka.



6. Check the components you want to install and uncheck the components you do not want to install. Click **Next** to continue.



7. 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.

8. Choose a Start Menu folder in which you would like to create the Weka 3.8.6 shortcuts. You can also enter a name to create a new folder. And then click on **Install** to continue.



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



10.**Installation Complete:** Setup was completed successfully. Click **Next** to continue.

11. Weka 3.8.6 has been installed on your computer. Click **Finish** to close Setup.



(b) Here are the features of some of its key components:

- **Explorer:** WEKA's graphical user interface (GUI) for data preprocessing, modelling, and evaluation, offering tools for visualization, classification, clustering, and more.
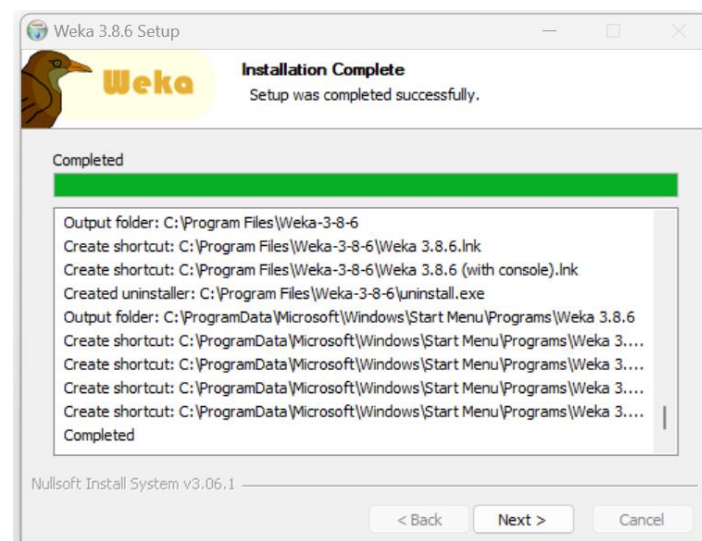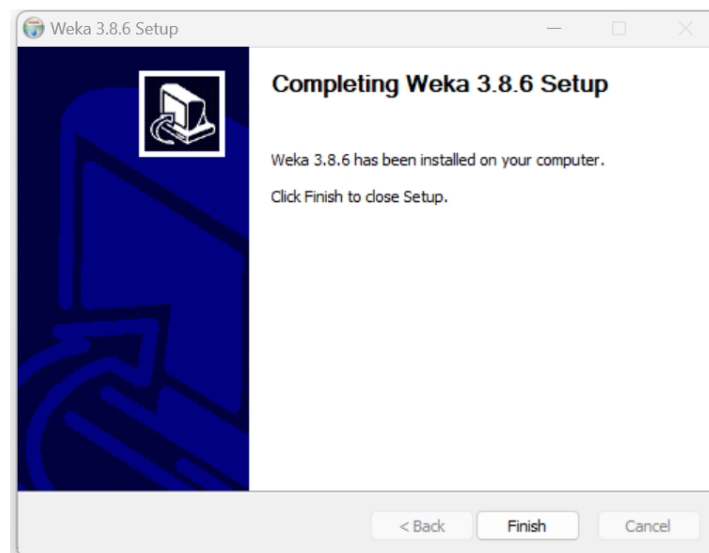- **Knowledge Flow Interface:** Enables users to design and execute machine learning workflows graphically, creating data processing pipelines with components like loaders, filters, classifiers, and evaluators.
- **Experimenter:** Facilitates systematic comparison and evaluation of machine learning algorithms across multiple datasets, automating experiment execution, result collection, and statistical analysis.
- **Command-Line Interface (CLI):** Allows users to perform WEKA operations via text-based commands, catering to scripting, batch processing, and integration into workflows for advanced users.



(c) Here's a brief overview of the options available in the WEKA Explorer for each panel:

- **Deselect Attributes Panel:** Allows users to manually select or deselect specific attributes from the dataset based on criteria like relevance or importance.
- **Preprocess Panel:** Enables various data preprocessing tasks such as handling missing values, normalization, discretization, and attribute selection to prepare the dataset for analysis.
- **Classify Panel:** Allows users to apply machine learning algorithms for classification tasks, selecting classifiers, setting parameters, and evaluating performance metrics.

- **Cluster Panel**: Enables the application of clustering algorithms to group similar instances in the dataset, with options to select algorithms, specify parameters, and visualize resulting clusters.
- **Associate Panel:** Used for association rule mining, discovering relationships between variables in large datasets, with options to specify parameters and view generated rules and metrics.
- **Visualize Panel:** Provides visualization tools for exploring and understanding the dataset and analysis results, including attribute distributions, scatter plots, decision boundaries, cluster assignments, and association rules.



(d) An Arff file contains two sections - header and data.

- The header describes the attribute types.
- The data section contains a comma separated list of data.

As an example for Arff format, the Weather data file loaded from the WEKA sample databases is shown below:



- The @relation tag defines the name of the database.
- The @attribute tag defines the attributes.
- The @data tag starts the list of data rows each containing the comma separated fields.

# PROGRAM 3:

To understand the working of datasets in WEKA and to perform demonstration of preprocessing on dataset weather.arff.

# WORKING:

Using the Open file ... option under the Preprocess tag select the weather-nominal.arff file.



When you open the file, your screen looks like as shown here –

Let us first look at the highlighted **Current relation** sub window. It shows the name of the database that is currently loaded. You can infer two points from this sub window −

- There are 14 instances - the number of rows in the table.

- The table contains 5 attributes - the fields, which are discussed in the upcoming sections.



On the left side, notice the **Attributes** sub window that displays the various fields in the database.



The **weather** database contains five fields - outlook, temperature, humidity, windy and play. When you select an attribute from this list by clicking on it, further details on the attribute itself are displayed on the right-hand side.

Let us select the temperature attribute first. When you click on it, you would see the following screen −



In the **Selected Attribute** sub window, you can observe the following −
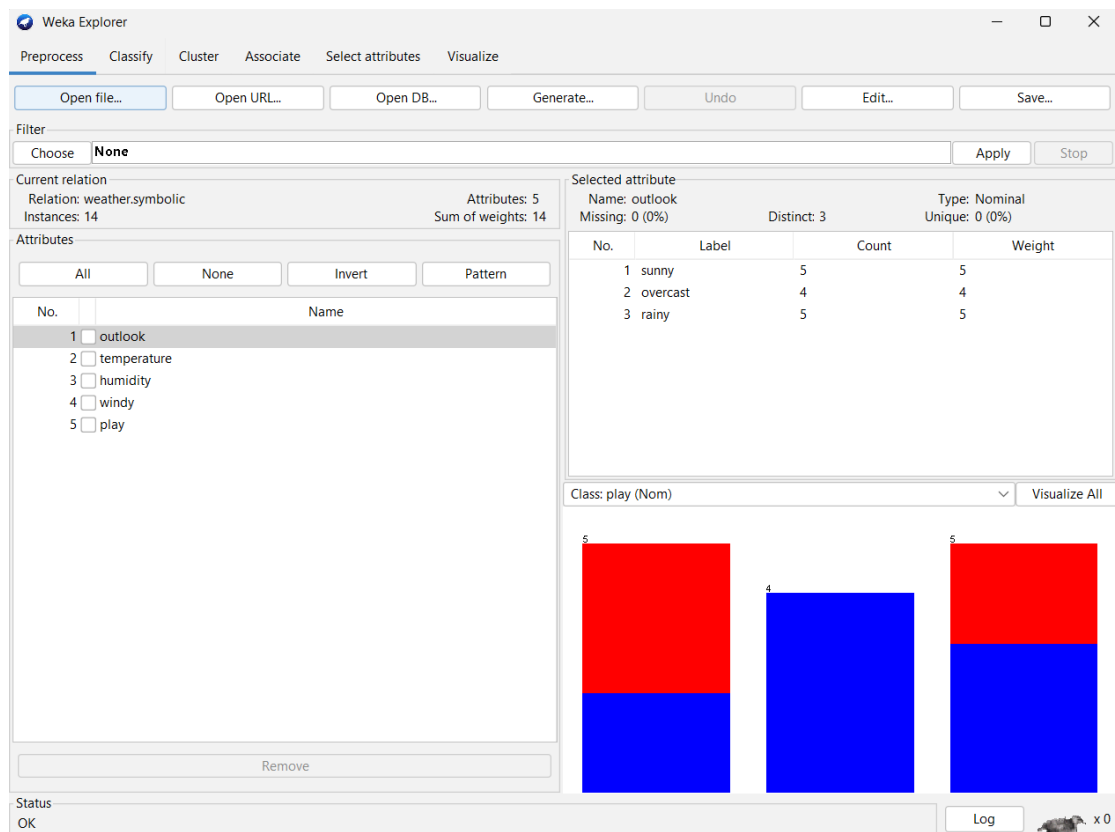
- The name and the type of the attribute are displayed.

- The type for the **temperature** attribute is **Nominal**.

- The number of **Missing** values is zero.

- There are three distinct values with no unique value.

- The table underneath this information shows the nominal values for this field as hot, mild, and cold.

- It also shows the count and weight in terms of a percentage for each nominal value.

At the bottom of the window, you see the visual representation of the **class** values.

If you click on the **Visualize All** button, you will be able to see all features in one single window as shown here −

All attributes

outlook | temperature | humidity | windy | play

# PROGRAM 4:

To apply Numeric Transform (data preprocessing step) on Iris Dataset.

# WORKING:

To apply Numeric Transform as a data preprocessing step on the Iris dataset in WEKA, you can follow these steps:

1. **Load the Dataset**: Open WEKA and navigate to the "Explorer" interface.

2. **Open the Iris Dataset**: Click on the "Open file..." button and select the Iris dataset file. The Iris dataset is a commonly used dataset available in the "data" folder of WEKA's installation directory.

3. **Explore the Dataset**: Switch to the "Preprocess" panel to view the dataset attributes and instances. Ensure that you have a clear understanding of the dataset's structure and contents before proceeding with preprocessing.



4. **Apply Numeric Transform**: To apply a Numeric Transform, follow these steps:

   - Select the "NumericTransform" filter from the "Filters" panel on the left-hand side.

   - Drag and drop the NumericTransform filter onto the preprocessing pipeline area in the middle of the interface.



   - Configure the NumericTransform filter by double-clicking on it in the preprocessing pipeline. Specify the transformation options you want to apply to the numeric attributes in the dataset.

   - Common transformations include standardization (mean normalization and scaling), normalization to a specific range, or any custom transformation you may require.

- Click on the "Apply" button to apply the Numeric Transform filter to the dataset.

# PROGRAM 5:

To understand the importance of CSV data and then load student academic record (CSV format) in Weka.

# THOERY:

CSV (comma-separated values) files are a common way to store tabular data, such as student academic records. They are easy to read and write and can be opened by a variety of software programs. CSV data is important because it is a standard format that can be used to exchange data between different programs. This makes it easy to share data with others, and to collaborate on projects. CSV data is also easy to store and archive, making it a good choice for long-term data storage.

# WORKING:

To load a student academic record (CSV format) in Weka, you can use the following steps:

1. Open the Weka Explorer, Click the "Open file" button, Select the CSV file that you want to load, and click the "Open" button.



2. The CSV file will be loaded into Weka, and you can view it in the "Preprocess" panel. You can then use the data in the CSV file to create a machine learning model, or to perform other data analysis tasks.

## PROGRAM 6:

To perform decision tree classification using J48 algorithm on weather.arff.

## THOERY:

**Decision trees** are also known as **Classification And Regression Trees (CART)**. They work by learning answers to a hierarchy of if/else questions leading to a decision. These questions form a tree-like structure, and hence the name.

Each node in the tree represents a question derived from the features present in your dataset. Your dataset is split based on these questions until the maximum depth of the tree is reached. The last node does not ask a question but represents which class the value belongs to.

- The topmost node in the Decision tree is called the **Root node.**

- The bottom-most node is called the **Leaf node.**

- A node divided into sub-nodes is called a **Parent node.** The sub-nodes are called **Child nodes.**

J48 is a decision tree algorithm implemented in WEKA, which is based on the C4.5 algorithm developed by Ross Quinlan.

## WORKING:

To perform decision tree classification using the J48 algorithm on the "weather.arff" dataset in WEKA, you can follow these steps:

1. Load the Dataset: Open WEKA and navigate to the "Explorer" interface.

2. Open the "weather.arff" Dataset: Click on the "Open file..." button and select the "weather.arff" dataset file.

3. Choose the Class Attribute: In the "Preprocess" panel, ensure that the class attribute (the attribute you want to predict) is correctly set. For the "weather.arff" dataset, the class attribute is likely to be something like "play," indicating whether or not to play outside based on weather conditions.

4. Select the Classify Panel: Switch to the "Classify" panel in WEKA.

5. Choose J48 Algorithm: In the "Classifier" section of the panel, click on the dropdown menu to select the J48 algorithm. J48 is WEKA's implementation of the C4.5 decision tree algorithm.

6. Set Options (if needed): Optionally, you can set specific options for the J48 algorithm by clicking on the "More options..." button. This allows you to customize parameters such as the confidence factor, minimum number of instances per leaf, and pruning options.

7. Run the Classifier: Once you have selected the J48 algorithm and set any desired options, click on the "Start" button to run the classifier. WEKA will train the decision tree classifier on the "weather.arff" dataset using the J48 algorithm.

8. Evaluate the Model: After the classifier has been trained, WEKA will display evaluation results in the "Classifier output" section. This includes metrics such as accuracy, precision, recall, and F-measure. You can use these metrics to assess the performance of the decision tree classifier on the dataset.

9. Visualize the Decision Tree (Optional): If you want to visualize the decision tree that was learned by the J48 algorithm, you can click on the "Visualize tree" button. This will display a graphical representation of the decision tree, allowing you to explore its structure and decision rules.

## PROGRAM 7:

Demonstration of classification rule process on dataset employee.arff using naïve bayes algorithm.

## THOERY:

**Naive Bayes** is a probabilistic classification algorithm based on Bayes' theorem, with an assumption of feature independence.

It calculates the probability of each class given a set of input features and selects the class with the highest probability as the predicted class. Despite its simplicity, Naive Bayes often performs well in practice, especially on text **classification** tasks and datasets with high-dimensional features.

It is computationally efficient and requires minimal parameter tuning, making it suitable for large datasets and real-time applications. However, its assumption of feature independence may not hold true in all datasets, potentially leading to suboptimal performance in some cases.

## WORKING:

To demonstrate the classification rule process using the Naïve Bayes algorithm on the "employee.arff" dataset in WEKA, follow these steps:

1. Load the Dataset: Open WEKA and navigate to the "Explorer" interface.

2. Open the "employee.arff" Dataset: Click on the "Open file..." button and select the "employee.arff" dataset file.

3. Choose the Class Attribute: In the "Preprocess" panel, ensure that the class attribute (the attribute you want to predict) is correctly set. For the "employee.arff" dataset, the class attribute might be something like "left," indicating whether an employee left the company or not.

4. Select the Classify Panel: Switch to the "Classify" panel in WEKA.

5. Choose Naïve Bayes Algorithm: In the "Classifier" section of the panel, click on the dropdown menu and select the Naïve Bayes algorithm. Naïve Bayes is a probabilistic classifier based on Bayes' theorem, assuming independence between features.

6. Run the Classifier: Once you have selected the Naïve Bayes algorithm, click on the "Start" button to run the classifier. WEKA will train the Naïve Bayes classifier on the "employee.arff" dataset.

7. Evaluate the Model: After the classifier has been trained, WEKA will display evaluation results in the "Classifier output" section. This includes metrics such as accuracy, precision, recall, and F-measure, which assess the performance of the Naïve Bayes classifier on the dataset.

8. Interpret the Classification Rules: While Naïve Bayes does not explicitly generate decision rules like some other algorithms, you can interpret the classification process based on probability estimates. Naïve Bayes calculates the probability of each class label given the input features and assigns the class label with the highest probability as the predicted class. You can analyze the conditional probabilities of features given each class to understand how the classifier makes decisions.

9. Visualize the Classifier Output (Optional): If you want to visualize the classification results or explore the probabilities assigned by the Naïve Bayes classifier, you can use visualization tools provided by WEKA.

# PROGRAM 8:

To understand the concept of discretization and to perform discretization on the dataset airline.arff.
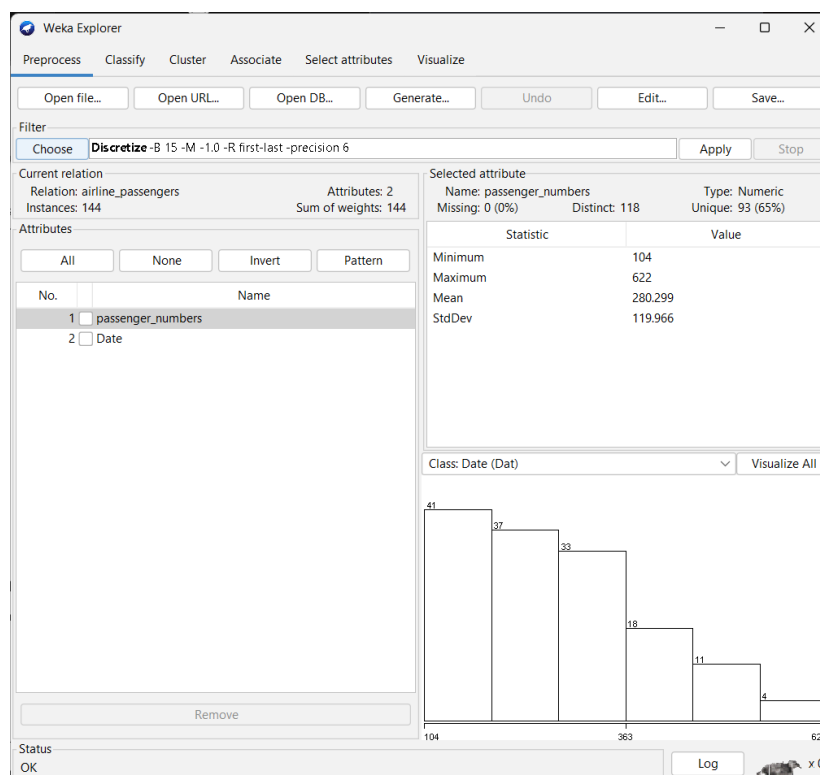
# THOERY:

Discretization is the process of converting continuous data into discrete intervals or categories. This is often done to simplify analysis or to prepare data for certain types of algorithms that require categorical input. There are several methods for discretization, including equal-width binning, equal-frequency binning, and clustering-based methods.

# WORKING:

Step-by-step guide to discretizing the dataset **airline.arff** in WEKA:

1. **Open WEKA**: Launch the WEKA application.

2. **Load Dataset**: Go to the "Explorer" tab, click on the "Open file" button, and select the **airline.arff** dataset.

3. **Preprocess Tab**: Switch to the "Preprocess" tab.

4. **Choose Filter**: In the Preprocess tab, you'll see a list of preprocessing filters on the left. Scroll down or type "discretize" in the search bar to find the "Discretize" filter.

5. **Apply Discretize Filter**: Drag the "Discretize" filter from the left panel to the right panel, where the preprocessing pipeline is constructed.

6. **Configure Discretize Filter (Optional)**: You can click on the "Discretize" filter in the right panel to configure its parameters. For example, you can specify the number of bins or select different discretization methods.



7. **Run**: Click on the "Start" button to run the preprocessing pipeline with the Discretize filter.

8. **Save Results**: After the preprocessing is complete, you can save the discretized dataset by clicking on the "Save" button.

# PROGRAM 9:

To create Training, Validation and Test dataset for iris.arff.

# THOERY:

Brief overview of training, validation, and test datasets and their roles in machine learning:

1. **Training Dataset**:

The training dataset is used to train machine learning models. It contains a labelled set of examples that the model learns from. The model adjusts its parameters based on the patterns and relationships present in the training data.

2. **Validation Dataset**:

The validation dataset is used to fine-tune the model's hyperparameters and assess its performance during training. It helps in preventing overfitting by providing an independent dataset to evaluate the model's generalization performance. After training on the training dataset, the model's performance is evaluated on the validation dataset, and hyperparameters are adjusted accordingly.

3. **Test Dataset**:

The test dataset is used to evaluate the final performance of the trained model. It provides an unbiased estimate of the model's performance on unseen data. The test dataset should be kept separate from the training and validation datasets to ensure an objective evaluation.

# WORKING:

Creating training, validation, and test datasets from the iris.arff dataset in WEKA involves splitting the dataset into three separate subsets. Here's how you can do it:
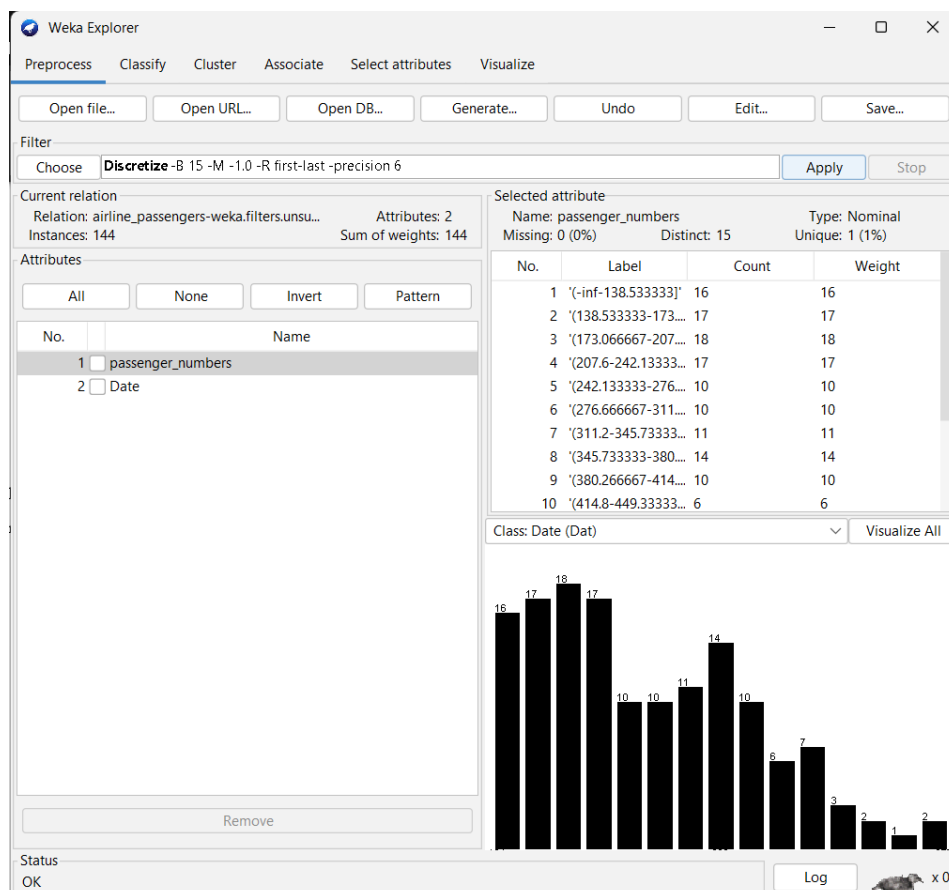
- Training Set:

1. Open WEKA: Launch the WEKA application.

2. Load Dataset: Go to the "Explorer" tab, click on the "Open file" button, and select the iris.arff dataset.

3. Select the **RemovePercentage** filter in the preprocess panel.
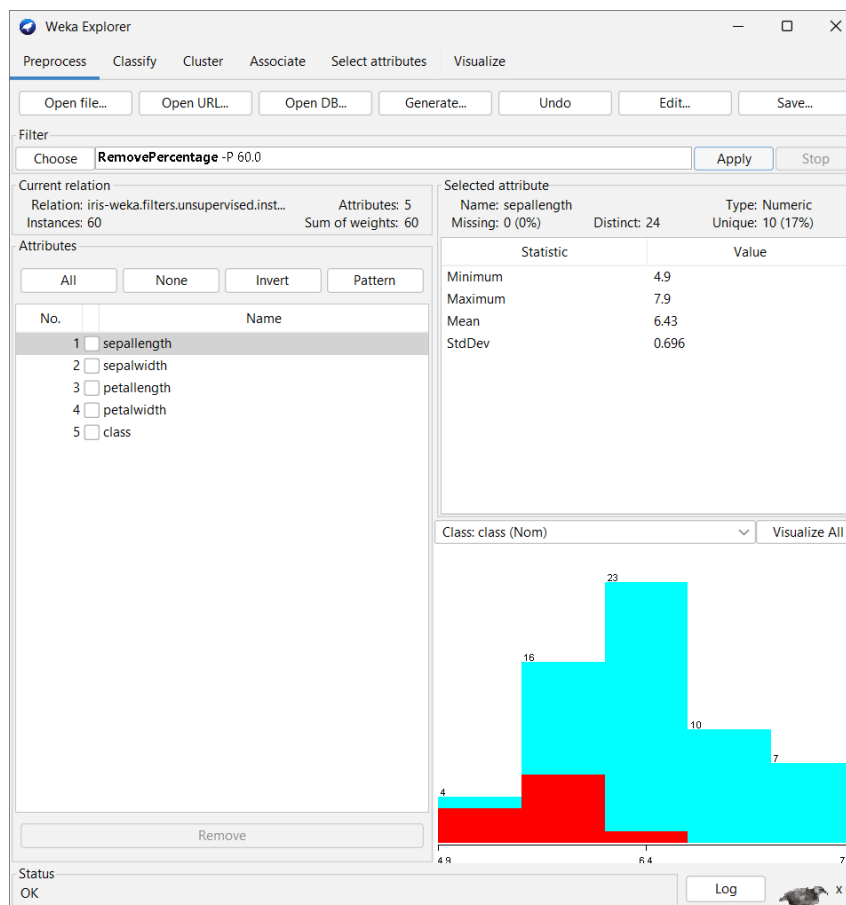4. Set the correct percentage for the split.
5. Apply the filter.
6. Save the generated data as a new file.



- Training Set:

1. Open WEKA: Launch the WEKA application.

2. Load Dataset: Go to the "Explorer" tab, click on the "Open file" button, and select the iris.arff dataset.
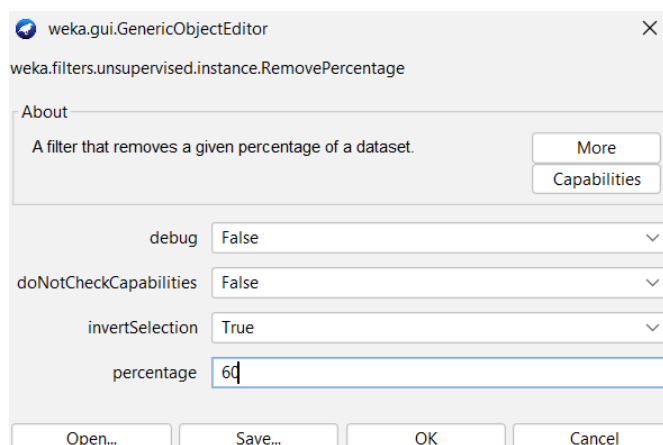
3. select the RemovePercentage filter if not yet selected.

4. set the invertSelection property to true.

5. apply the filter.

6. save the generated data as new file.

# PROGRAM 10:

To apply Apriori technique on the dataset and to generate association rules.

# THOERY:

The Apriori algorithm is a common method for association rule mining. It works by generating all possible itemsets in a dataset, finding the frequent itemsets using the "join and prune" technique, and generating association rules from the frequent itemsets.

The Apriori algorithm is a powerful tool for association rule mining. It can be used to find interesting and useful patterns in data. However, it can also be computationally expensive, especially for large datasets.

# WORKING:

To apply the Apriori algorithm and generate association rules in WEKA, follow these steps:

1. **Load the Dataset**: Open WEKA and load your dataset.

2. **Explorer Tab**: Go to the "Explorer" tab.

3. **Choose Apriori Algorithm**: In the "Choose" section, click on the "Associate" button.

4. **Select Apriori**: From the list of association algorithms, select "Apriori" by clicking on it. This will open the configuration panel for the Apriori algorithm.

5. **Configure Apriori Parameters**:

   - You can adjust various parameters such as minimum support, minimum confidence, and other settings based on your dataset and requirements.

   - Set the minimum support and minimum confidence thresholds according to your desired levels of significance.

6. **Run Apriori**: Click on the "Start" button to run the Apriori algorithm on your dataset.

7. **View Association Rules**: Once the algorithm finishes running, you'll see the generated association rules in the "Association rules" tab below the configuration panel.

# PROGRAM 11:

(a) Demonstration of classification rule process on dataset student.arff using J48 algorithm.

(b) Demonstration of classification rule process on dataset employee.arff using J48 algorithm.
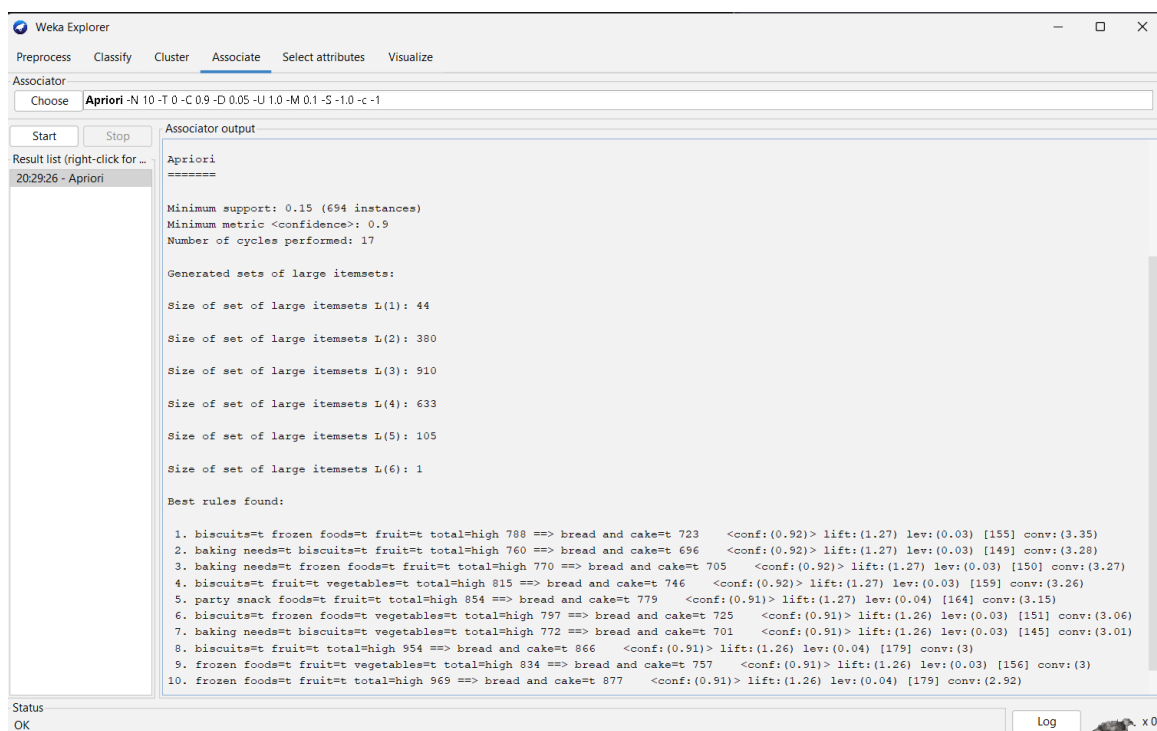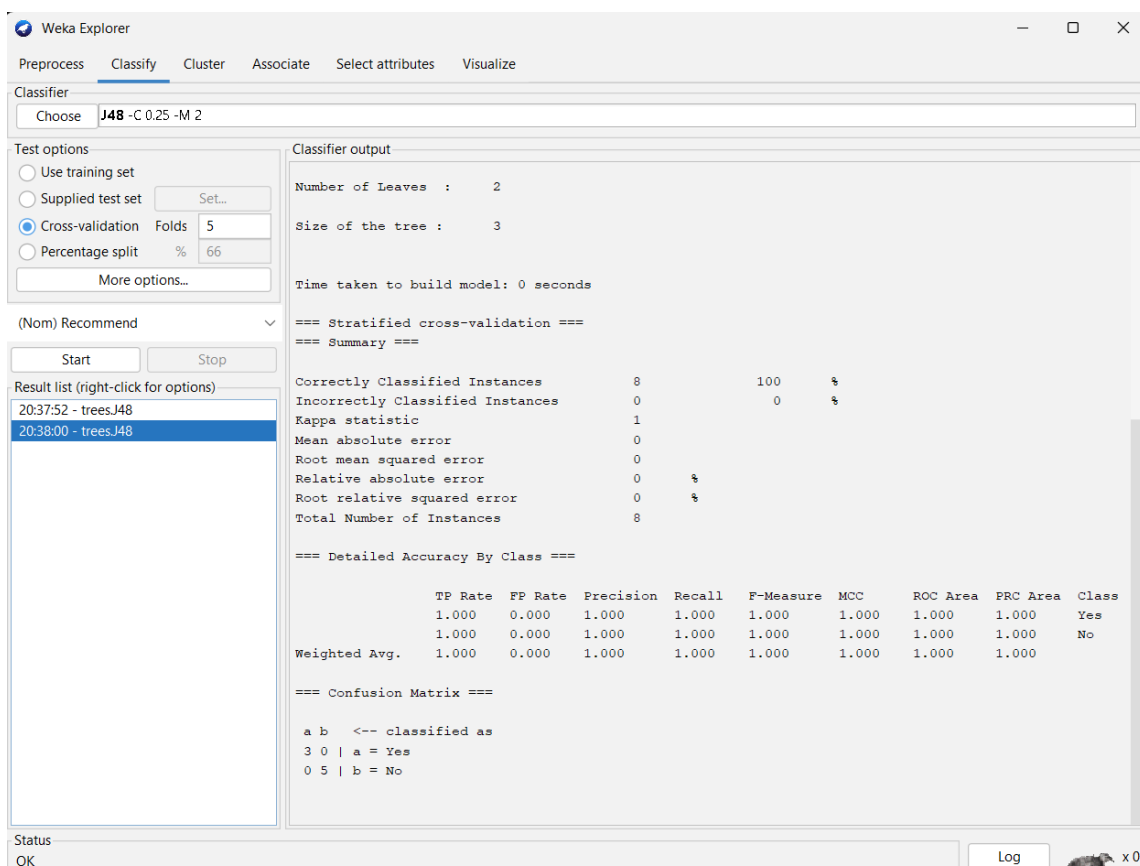
## THOERY:

The Apriori algorithm is a common method for association rule mining. It works by generating all possible itemsets in a dataset, finding the frequent itemsets using the "join and prune" technique, and generating association rules from the frequent itemsets.

The Apriori algorithm is a powerful tool for association rule mining. It can be used to find interesting and useful patterns in data. However, it can also be computationally expensive, especially for large datasets.

## WORKING:

(a) To demonstrate the classification rule process using the J48 algorithm on the dataset **student.arff** in WEKA, follow these steps:

1. **Load the Dataset**: Open WEKA and load the **student.arff** dataset.

2. **Explorer Tab**: Go to the "Explorer" tab.

3. **Choose J48 Algorithm**: In the "Choose" section, click on the "Classify" button.

4. **Select J48**: From the list of classification algorithms, select "J48" by clicking on it. This will open the configuration panel for the J48 algorithm.

5. **Configure J48 Parameters**:

   - You can adjust various parameters such as confidence factor, minimum number of instances per leaf, and other settings based on your dataset and requirements.

6. **Run J48**: Click on the "Start" button to run the J48 algorithm on your dataset.

(b) To demonstrate the classification rule process using the J48 algorithm on the dataset **employee.arff** in WEKA, you can follow a similar process as described earlier. Here's a step-by-step guide:

1. **Load the Dataset**: Open WEKA and load the **employee.arff** dataset.

2. **Explorer Tab**: Go to the "Explorer" tab.

3. **Choose J48 Algorithm**: In the "Choose" section, click on the "Classify" button.

4. **Select J48**: From the list of classification algorithms, select "J48" by clicking on it. This will open the configuration panel for the J48 algorithm.

5. **Configure J48 Parameters**:

    - Adjust various parameters such as confidence factor, minimum number of instances per leaf, and other settings based on your dataset and requirements.

    - Set the parameters according to your desired model complexity and performance requirements.

6. **Run J48**: Click on the "Start" button to run the J48 algorithm on your dataset.

---

**Weka Explorer**  — □ ✕

Preprocess  Classify  Cluster  Associate  Select attributes  Visualize

**Classifier**

Choose | J48 -C 0.25 -M 2

**Test options**
- ○ Use training set
- ○ Supplied test set    Set...
- ● Cross-validation   Folds  10
- ○ Percentage split    %  66

More options...

(Nom) class  ⌄

Start | Stop

Result list (right-click for options)
21:00:54 - trees.J48

**Classifier output**

```
Number of Leaves  :     3

Size of the tree :     5


Time taken to build model: 0 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances          42               73.6842 %
Incorrectly Classified Instances        15               26.3158 %
Kappa statistic                          0.4415
Mean absolute error                      0.3192
Root mean squared error                  0.4669
Relative absolute error                 69.7715 %
Root relative squared error             97.7888 %
Total Number of Instances               57

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC     ROC Area  PRC Area  Class
                 0.700    0.243    0.609      0.700   0.651      0.444   0.695     0.559     bad
                 0.757    0.300    0.824      0.757   0.789      0.444   0.695     0.738     good
Weighted Avg.    0.737    0.280    0.748      0.737   0.740      0.444   0.695     0.675

=== Confusion Matrix ===

  a  b   <-- classified as
 14  6 |   a = bad
  9 28 |   b = good
```

**Status**
OK                                                    Log     x 0