

Name Harshita Sharma
Sub. Data Mining
BRANCH- AIDS-B
Enroll no.: 21117711922 School _____



VIVEKANANDA INSTITUTE OF PROFESSIONAL STUDIES - TECHNICAL CAMPUS

Grade **A++** Accredited Institution by NAAC

NBA Accredited for MCA Programme; Recognized under Section 2(f) by UGC;
Affiliated to GGSIP University, Delhi; Recognized by Bar Council of India and AICTE
An ISO 9001:2015 Certified Institution

SCHOOL OF ENGINEERING & TECHNOLOGY

BTECH Programme: AIDS-B

Course Title: Data Mining

Course Code: AIDS305

Submitted To

Dr. Lakshita Aggarwal

Submitted By

**Name: Harshita Sharma
Enrollment No:21117711922**



VIVEKANANDA INSTITUTE OF PROFESSIONAL STUDIES - TECHNICAL CAMPUS

Grade A++ Accredited Institution by NAAC

NBA Accredited for MCA Programme; Recognized under Section 2(f) by UGC;
Affiliated to GGSIP University, Delhi; Recognized by Bar Council of India and AICTE
An ISO 9001:2015 Certified Institution

SCHOOL OF ENGINEERING & TECHNOLOGY

BTECH Programme: AIDS-B

Course Title: Data Mining

Course Code: AIDS305

Submitted To

Dr. Lakshita Aggarwal

Submitted By

**Name: Harshita Sharma
Enrollment No:21117711922**



VIVEKANANDA INSTITUTE OF PROFESSIONAL STUDIES - TECHNICAL CAMPUS
Grade A++ Accredited Institution by NAAC
NBA Accredited for MCA Programme; Recognized under Section 2(f) by UGC;
Affiliated to GGSIP University, Delhi; Recognized by Bar Council of India and AICTE
An ISO 9001:2015 Certified Institution
SCHOOL OF ENGINEERING & TECHNOLOGY

VISION OF INSTITUTE

To be an educational institute that empowers the field of engineering to build a sustainable future by providing quality education with innovative practices that supports people, planet and profit.

MISSION OF INSTITUTE

To groom the future engineers by providing value-based education and awakening students' curiosity, nurturing creativity and building capabilities to enable them to make significant contributions to the world.



VIVEKANANDA INSTITUTE OF PROFESSIONAL STUDIES - TECHNICAL CAMPUS
Grade A++ Accredited Institution by NAAC
NBA Accredited for MCA Programme; Recognized under Section 2(f) by UGC,
Affiliated to GGSIP University, Delhi; Recognized by Bar Council of India and AICTE
An ISO 9001:2015 Certified Institution
SCHOOL OF ENGINEERING & TECHNOLOGY

INDEX

Experiment Name	Date	Marks			Remark	Updated Marks	Faculty Signature
		Laboratory Assessment (15 Marks)	Class Participation (5 Marks)	Viva (5 Marks)			
Install & configuration of WEKA	5/8/24	15	5	5			Accepted 22/8/24
Preprocessing, mining, integration transformation on ARFF files using WEKA	12/8/24	13	4	4			Accepted 19/8/24
Apply association rule mining on dataset using WEKA	12/8/24	14	4	4			Accepted 20/8/24
Apply visualisation in data using WEKA	9/9/24	15	4	5			Accepted 20/9/24
Apply clustering on dataset using WEKA	9/9/24	14	4	4			Accepted 20/9/24
Apply classification in dataset, WEKA	23/9/24	15	4	4			Accepted 23/9/24
Evaluate the performance of classification techniques using diff parameters	23/9/24	14	4	5			Accepted 23/9/24



VIVEKANANDA INSTITUTE OF PROFESSIONAL STUDIES - TECHNICAL CAMPUS
Grade A++ Accredited Institution by NAAC
NBA Accredited for MCA Program
Affiliated to Anna University

NBA Accredited for MCA Programme; Recognized under Section 2(f) & 12(b) of UGC Act, 1956.

Admitted to GGSIP University, Delhi; Recognized by Bar Council of India and AICTE
An ISO 9001:2015 Certified Institute

An ISO 9001:2015 Certified Institution

SCHOOL OF ENGINEERING & TECHNOLOGY

INDEX

EXPERIMENT-1

AIM: Introduction & Installation of WEKA tools.

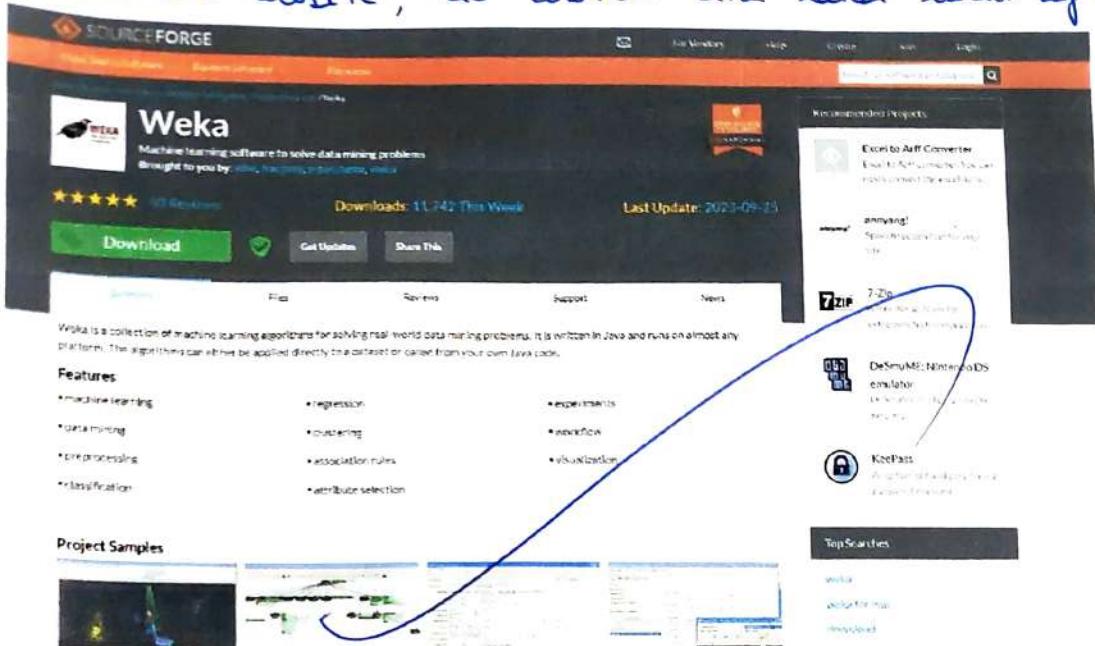
Theory: Weka environment for knowledge analysis (WEKA)
 It is a collection of machine learning & data ~~science~~ analysis free software licensed under the GNU General Public License.
 In other words, WEKA is a collection of ML algorithms for data mining tasks.
 It contains tools for data pre-processing, classification, regression, clustering, association rules & visualisation.

Weka provides the implementation of several algorithms one would select an algorithm of our choice, set the desired parameters and run it on dataset.

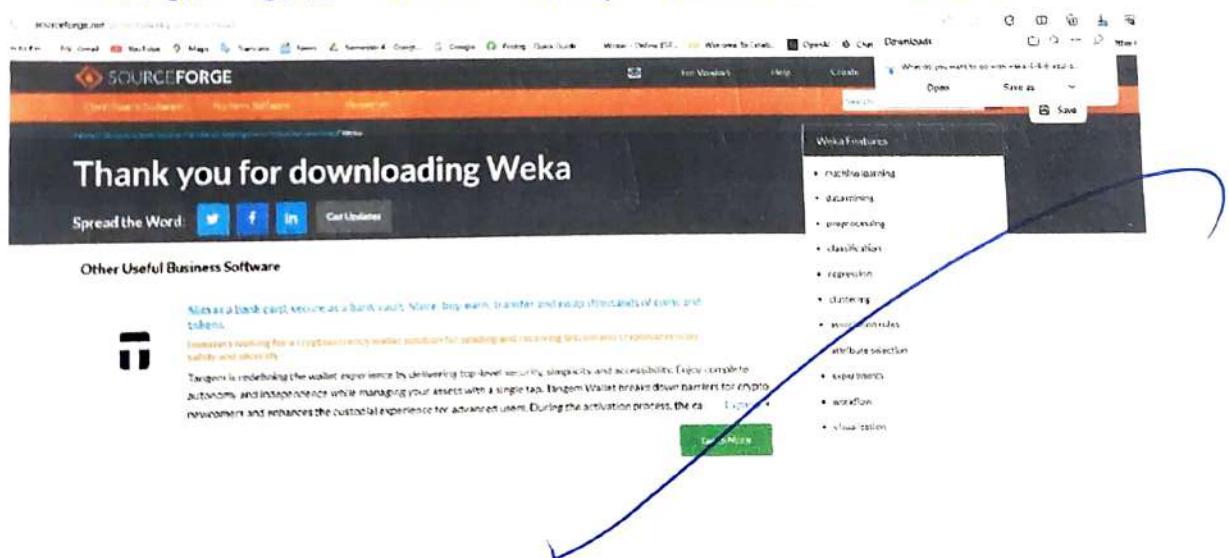
Thus, the use of WEKA results in a quick development of machine learning algorithms on the whole.

→ WEKA offers a graphical user interface (GUI) that allows users to interact with the software without needing to write the code, making it accessible to the non-programmers.

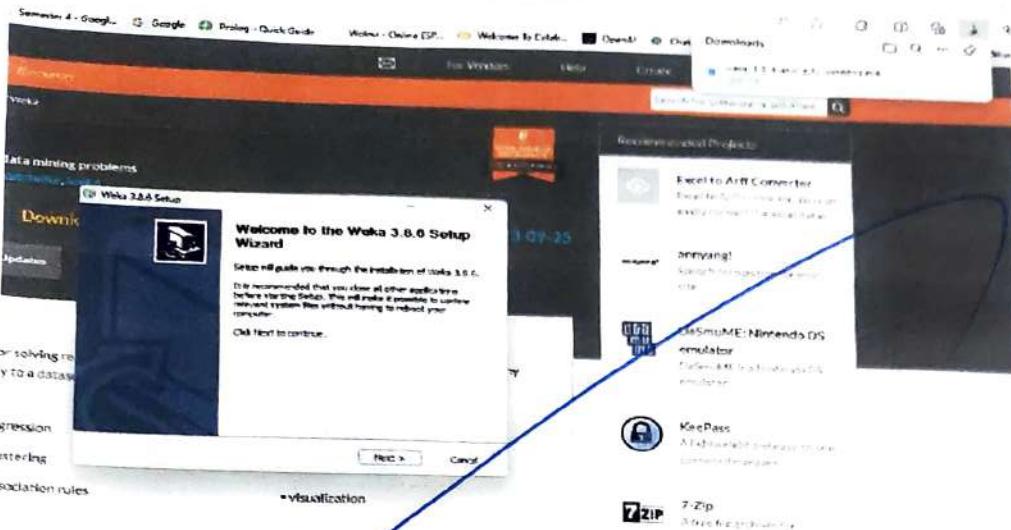
- Visit the website, to install the latest version of WEKA.



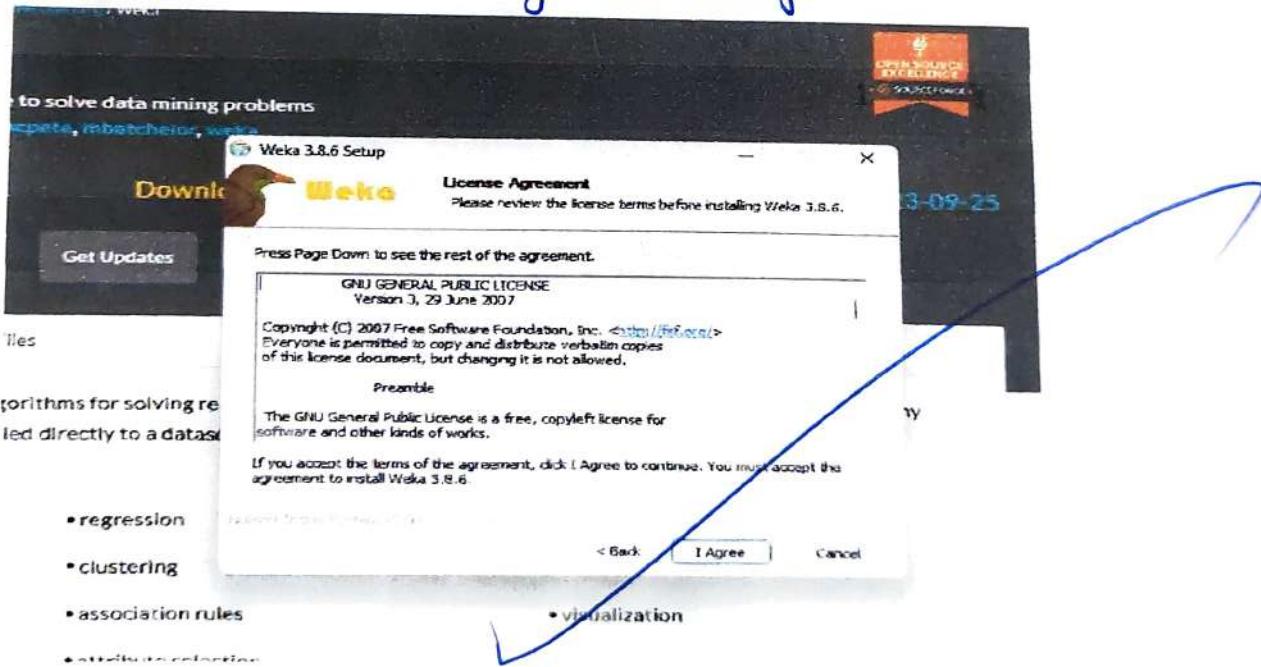
- Initiate the weka setup installation wizard.



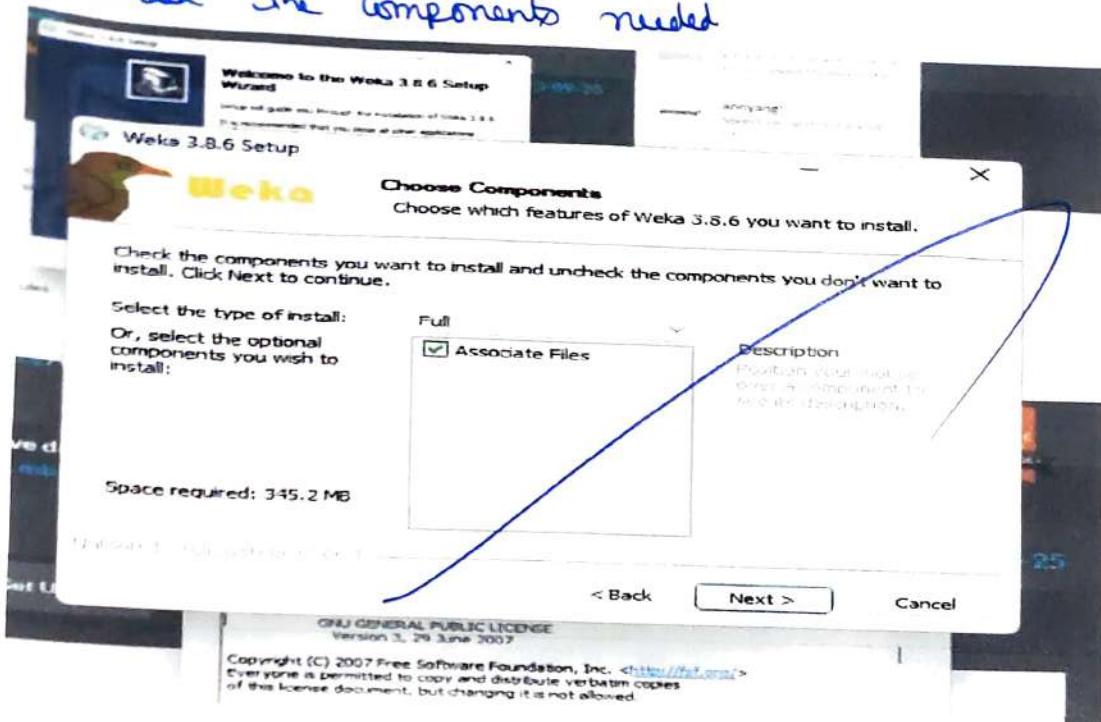
3. Accept the necessary terms & conditions



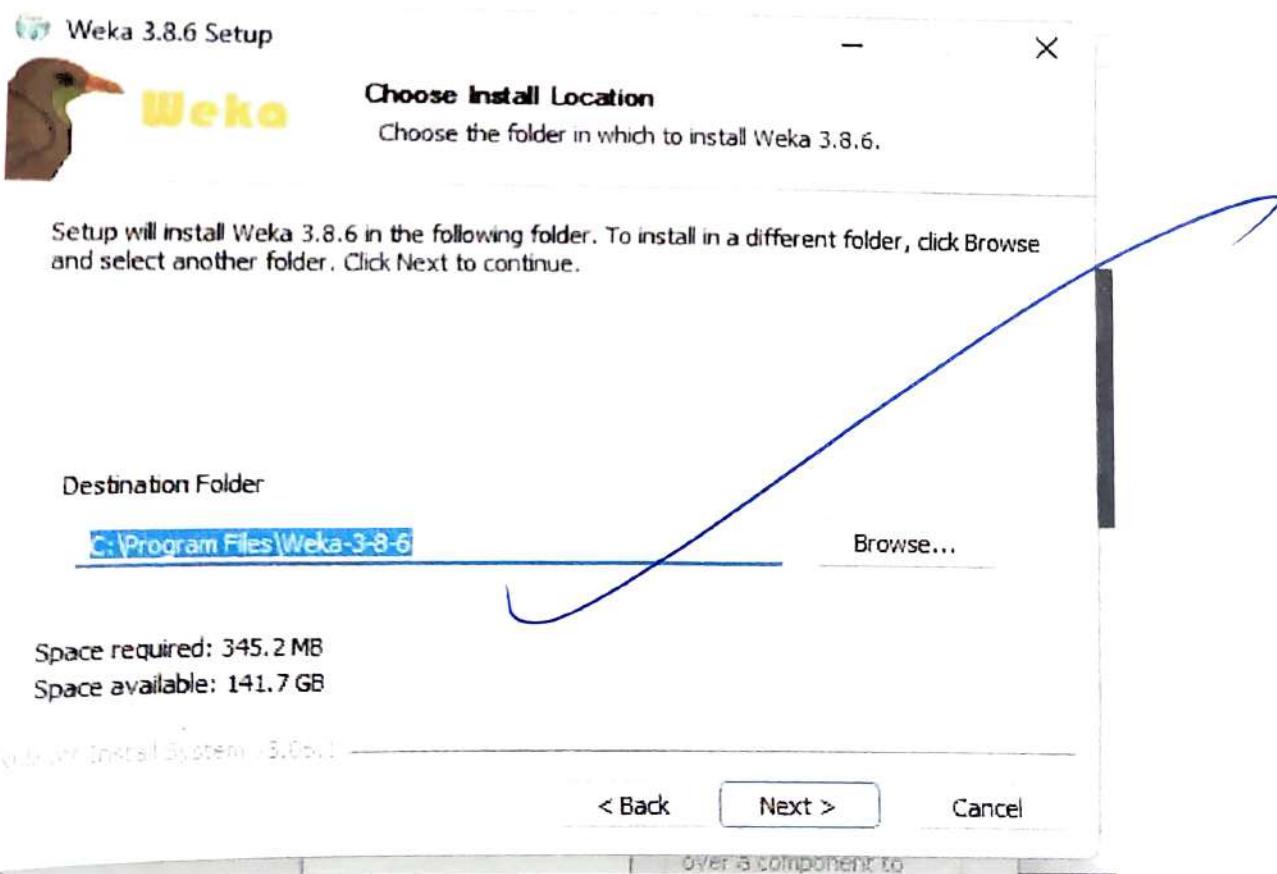
4. Setup weka storage directory.



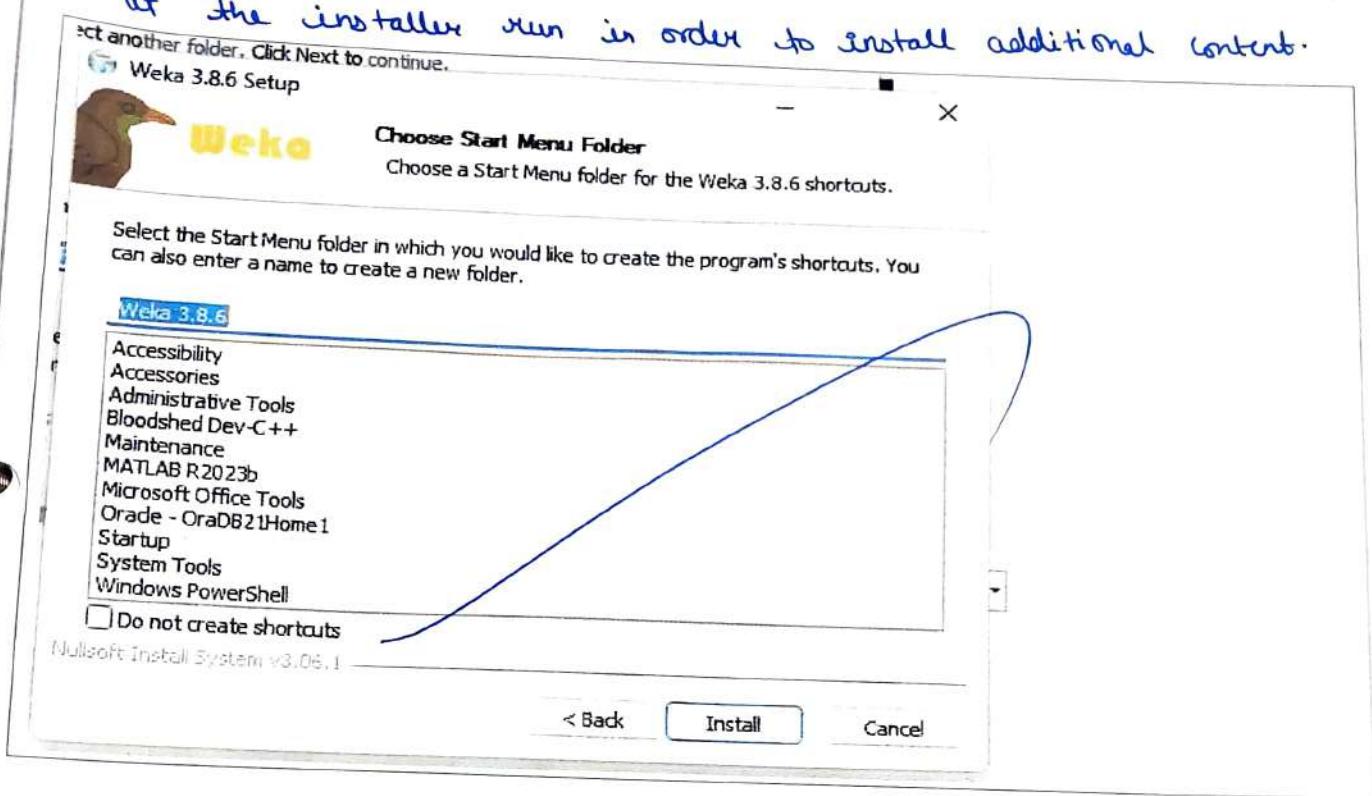
5. Choose the components needed



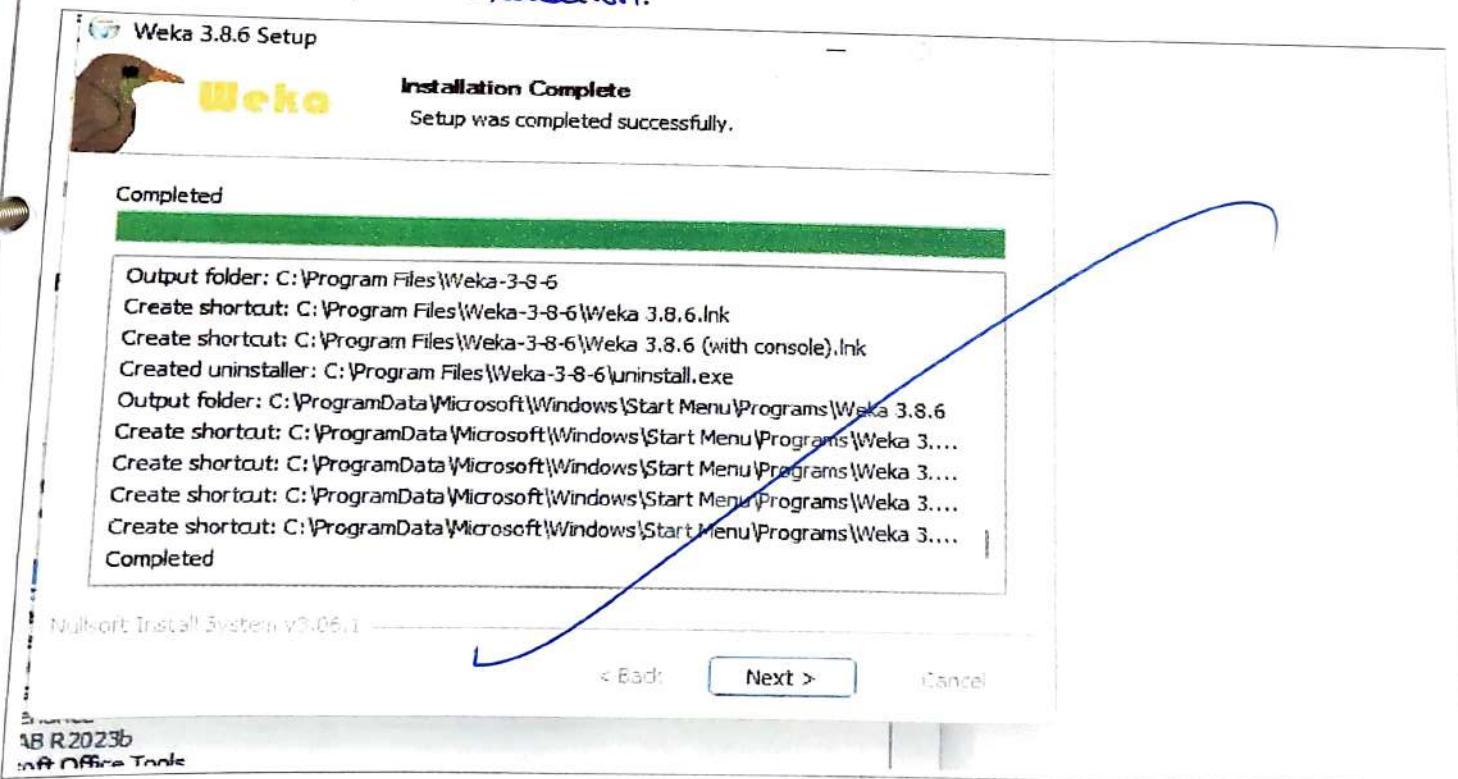
6. Choose the installation location



7. Let the installer run in order to install additional content.

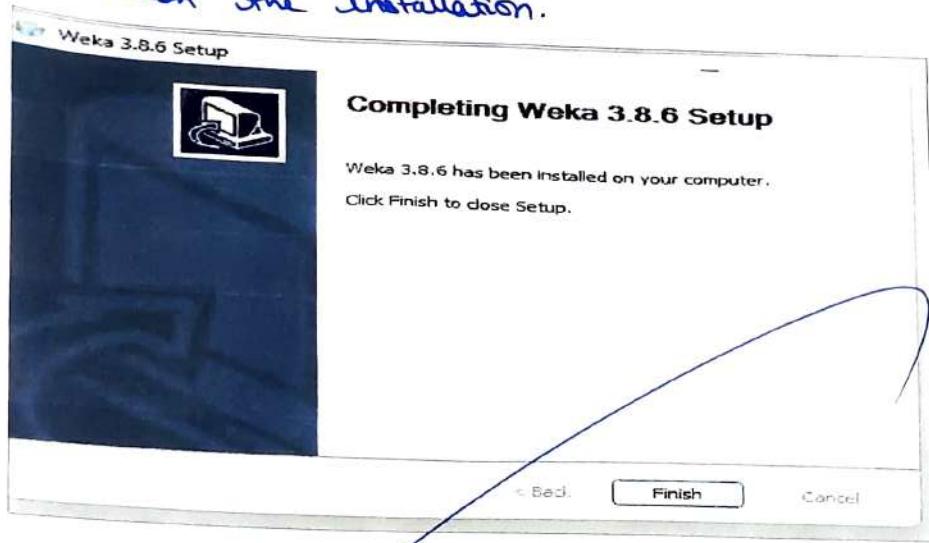


8. finish the installation.



9. Finish the installation.

6



10. WEKA is installed.



Learning Outcomes:

- successfully installed + understood the installation procedure of WEKA.
- learnt about the functions of WEKA on machine learning.

Weka
after

Experiment 2

Aim: Perform data pre-processing including cleaning, integration and transformation on ARFF files using WEKA.

Theory:

Data Preprocessing: It is the process of preparing raw data for analysis by handling missing values, outliers & noise to ensure high-quality input for models. It also helps improve accuracy and reliability of data set.

Data Cleaning: It is the process of correcting or deleting inaccurate, damaged, improperly formatted, duplicated or insufficient data from a dataset. Data cleaning lowers errors & raised caliber of the data. It is crucial to construct of a model.

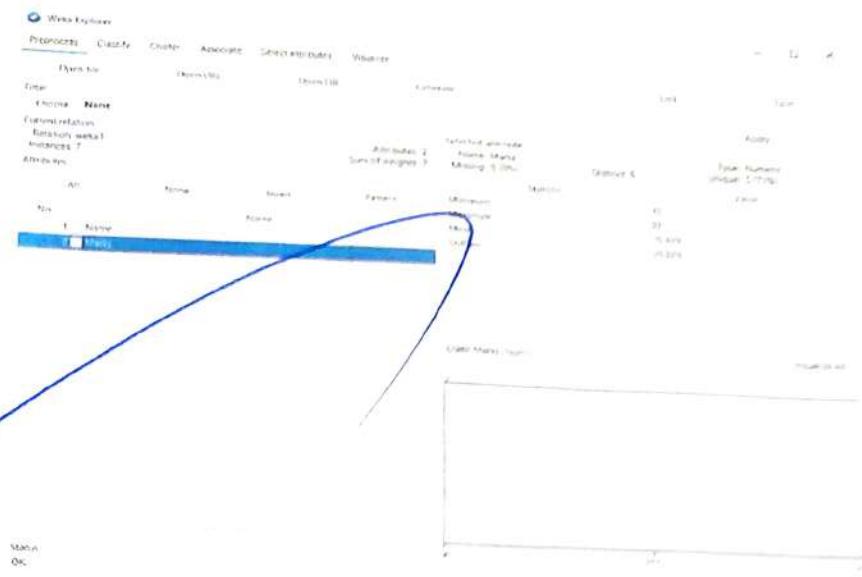
Data Integration: It is the process of merging data from several dispersed sources. Data integration is a second preprocessing method that includes merging data from the couple of heterogeneous data sources into coherent data to retain and provide a unified perspective of data.

Data Reduction: Data reduction techniques ensure the integrity of data while reducing the data. Data reduction is a process that reduces the volume of original data & represents it in a much smaller volume. Data reduction are used to obtain a reduced representation of dataset.

Data Transformation: Converts data into a suitable format or structure for analysis including normalisation, scaling and feature engineering. This process ensures the data to compatible with machine learning models.

	Name	Marks
1	Ram	89
2	shyam	90
3	lakhan	92
4	makan	67
5	lakhshman	45
6	shivam	43
7	rohan	67

.arff File



.text File

```

File Edit View

@relation weka2

@attribute Name string
@attribute Marks numeric

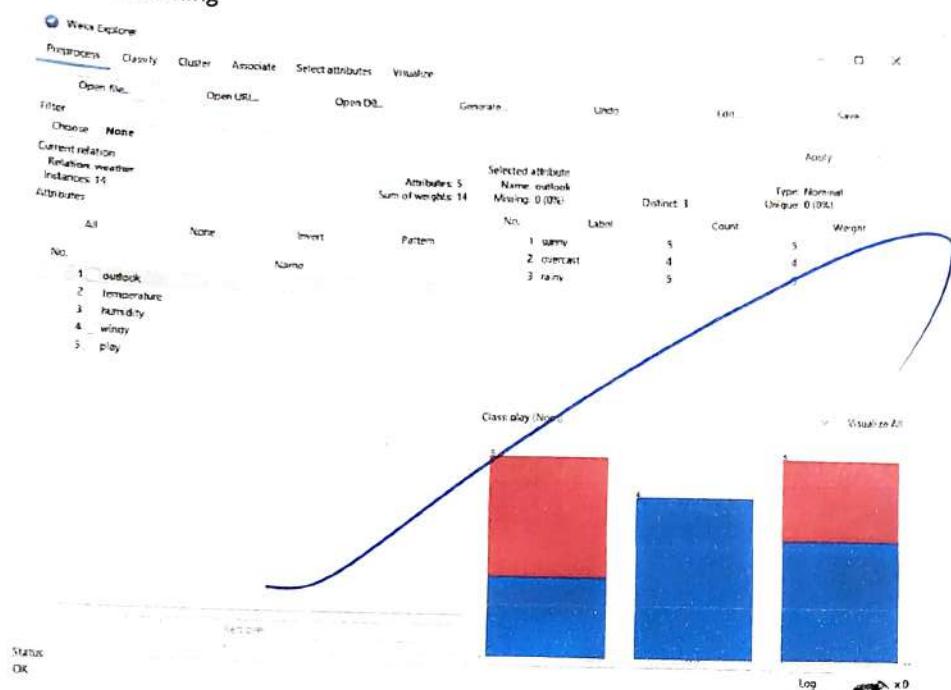
@data
Sita,67
Ram,78
    
```

.arff File

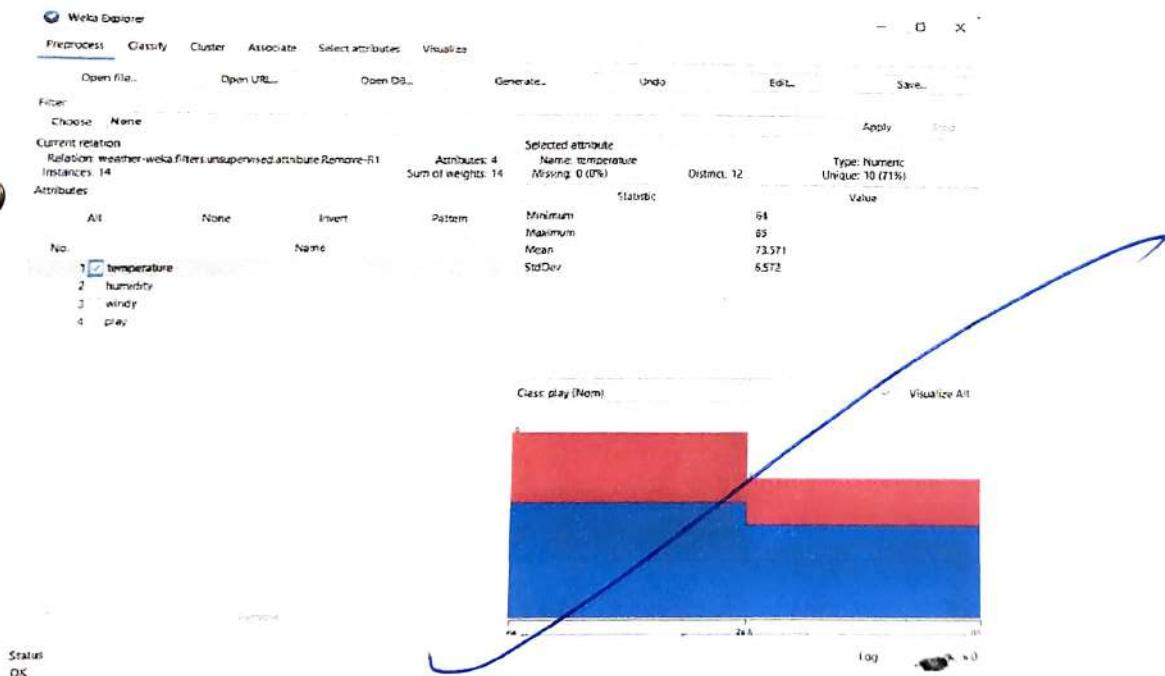


Data Cleaning

Before Cleaning

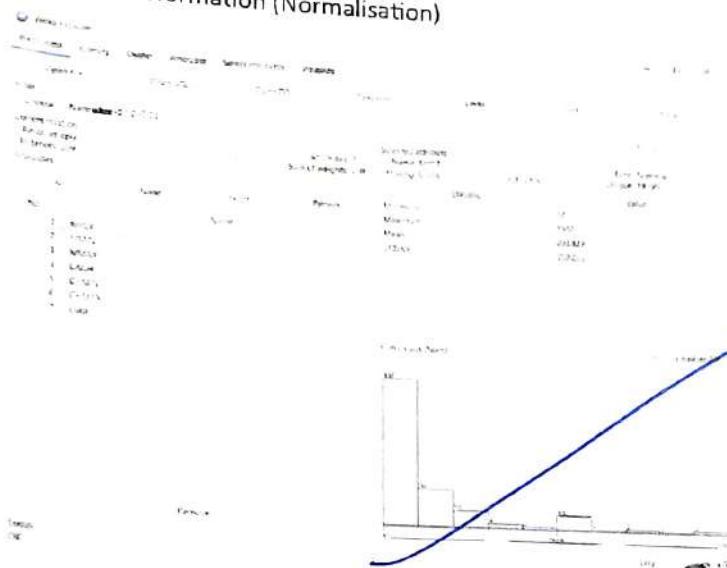


After Cleaning

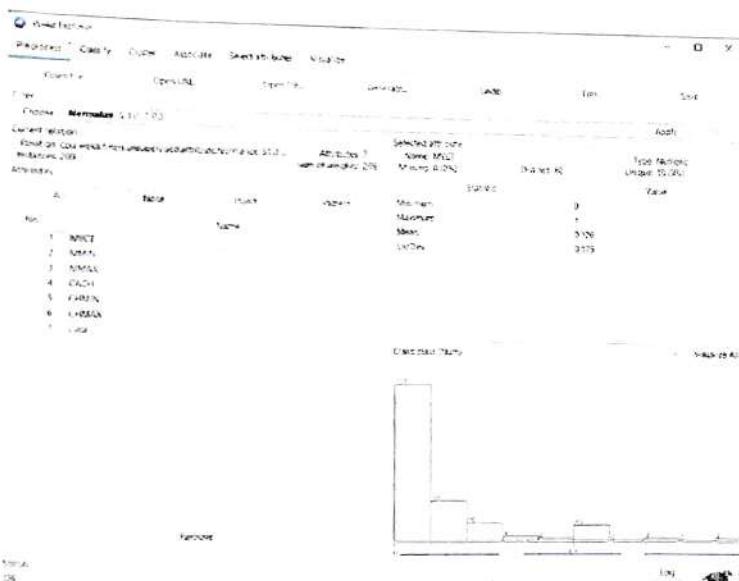


Data Transformation

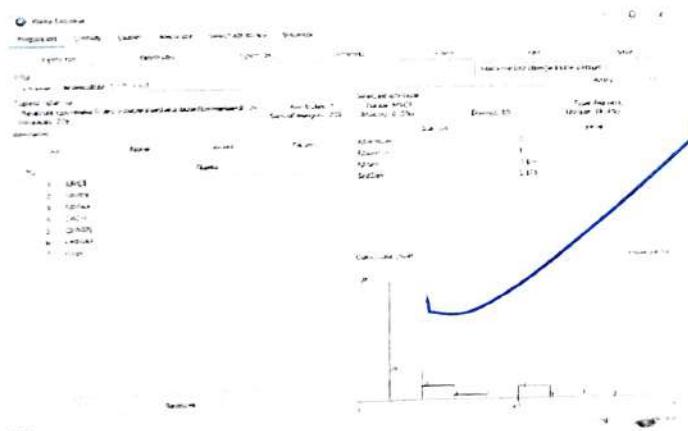
Before Transformation (Normalisation)



After Transformation



Undo



Data Integration

File 1 (Weka1.arff)



File 2 (Weka2.arff)



```
SimpleCLI  
Welcome to the WEKA SimpleCLI.  
Enter commands in the textfield at the bottom of  
the window. Use the up and down arrows to move  
through previous commands.  
Command completion for classnames and files is  
initiated with tab. In order to distinguish  
between files and classnames, file names must  
be either absolute or start with ./ or ..  
The latter is a shorthand for the home directory.  
Tab-space-space is used for deleting the text  
in the commandline in chunks.  
Type 'help' followed by Enter to see an overview  
of all commands.  
java weka.gui.FastEstimator arfffile testfile optionalargs -> classifier optionalargs  
Finished running input to classifier like weka3.arff
```

Learning Outcome:

- learnt about applying various machine learning techniques on weka such as data cleaning, data reduction, integration, etc.
- learnt converting .xls file onto .arff file + .tt file to .arff file

Harshita
9/9/2024

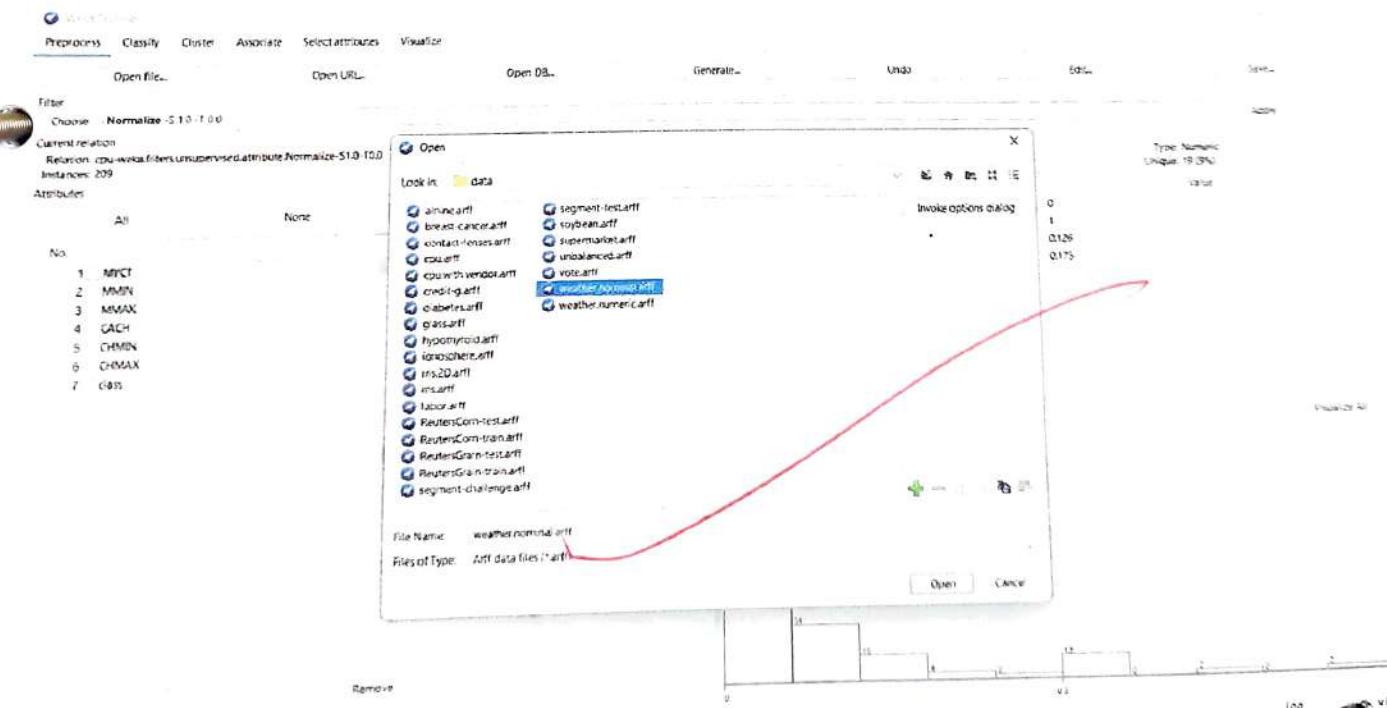
Experiment 3

Aim: Apply Association Rule Mining on a dataset using Weka.

Theory: Association Rule, is a data mining technique used to identify relationships between variables in large datasets.

- It uncovers patterns in the form of 'if-then' statements, showing how the occurrence of one item is associated with another.
- These rules are commonly used in market basket analysis to discover product associations, like 'if a customer buys bread, they are likely to buy butter'.
- support: measures how frequently an itemset appears in dataset.
- confidence: indicates the likelihood that an item B will be bought if item A is bought.
- lift: - shows how much more likely item B is to be bought with item A.

Opening Data:



Applying Apriori Algorithm

Week	Learning	Resources	Skills	Output	Assessments	Interventions	Feedback
Week 1	Introduction to Python	Python Tutorials	Basic Syntax	Code Examples	Quizzes	Online Help	Peer Review
Week 2	Control Flow	Control Flow Examples	Loops, Conditionals	Script Output	Code Review	Code Editor Plugins	Code Comments
Week 3	Data Structures	Data Structure Examples	Lists, Dictionaries, Sets	Program Output	Unit Tests	IDE Integration	API Documentation
Week 4	Object-Oriented Programming	Object-Oriented Examples	Classes, Inheritance	Code Snippets	Code Review	Code Editor Plugins	Code Comments
Week 5	Advanced Topics	Advanced Topics Examples	Decorators, Generators	Program Output	Code Review	Code Editor Plugins	Code Comments

Changing num Rules

New Rules

3

Learning Outcome:

- Understood the concept of Association Rule.
 - Learned how to apply various association Rules on WEKA.
 - Learned how to set parameters of an algorithm on WEKA.

Exhibit no 92a

Experiment 4

Aim: Apply Visualization on a dataset using Weka.

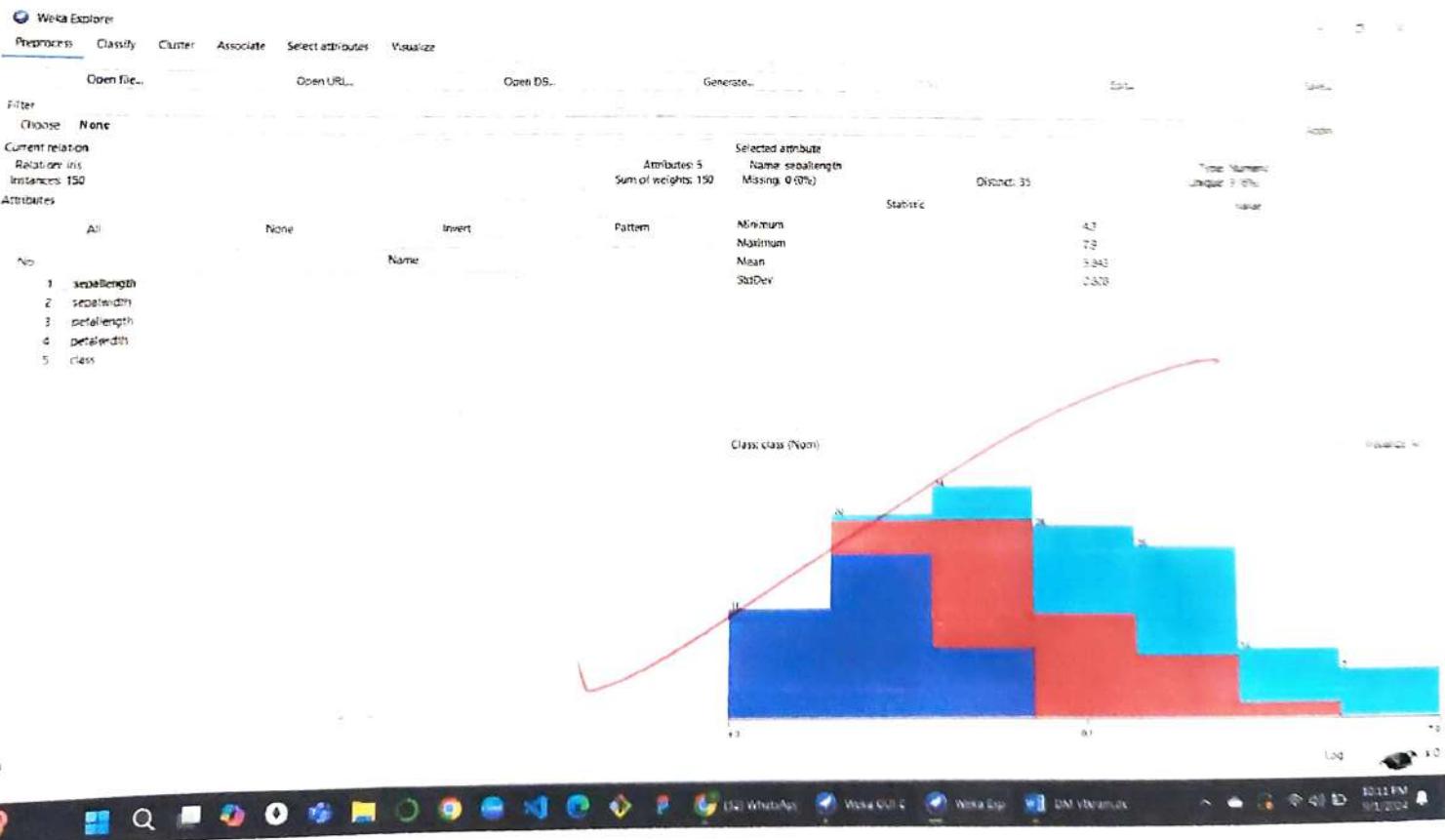
Theory: Data visualization is an important feature for ~~explaining~~ explaining and understanding datasets before applying ML algorithms. WEKA provides several tools and methods for visualizing the data, which can help in tasks such as data preprocessing, feature selection, and model evaluation.

* types of visualization in WEKA:

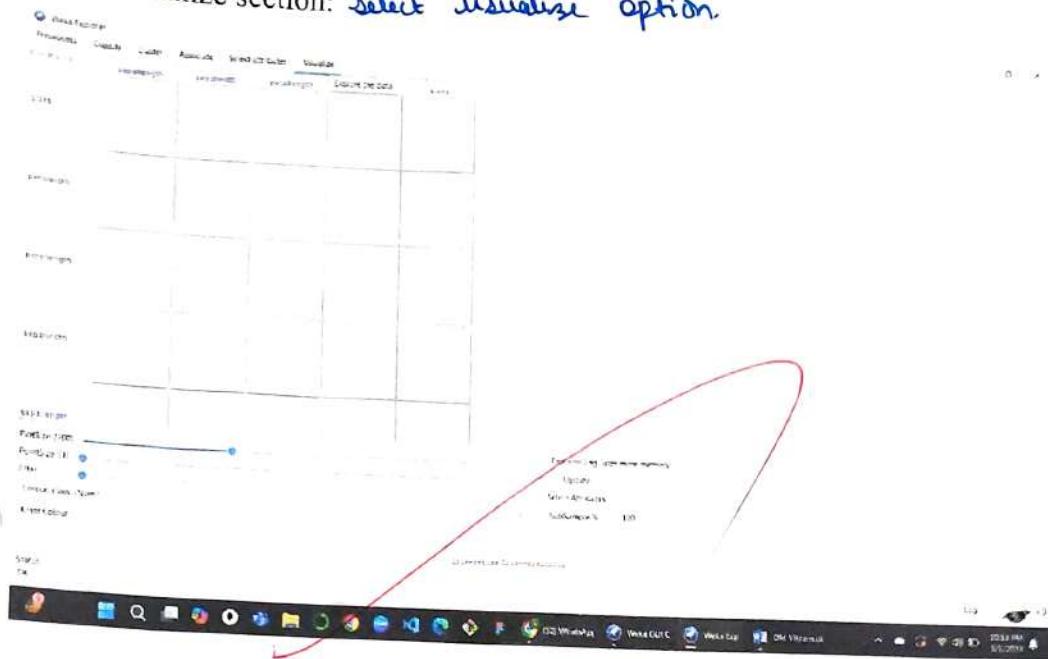
1. Data Explore: Histogram, Box Plot, ScatterPlot.
2. Visualisation Tools: Attribute Selection, Class distribution.
3. Model evaluation: ROC curves, confusion matrix.

Open Iris Dataset:

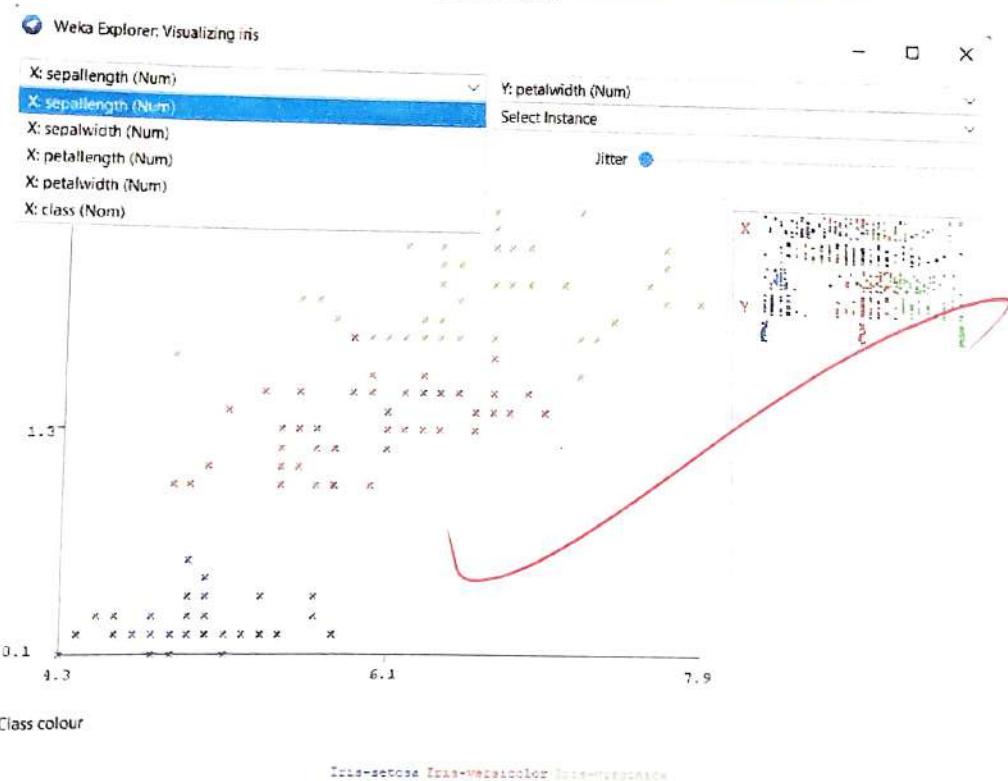
choose any numeric dataset from program files > weka > data > "any data"



Open visualize section: select visualize option.

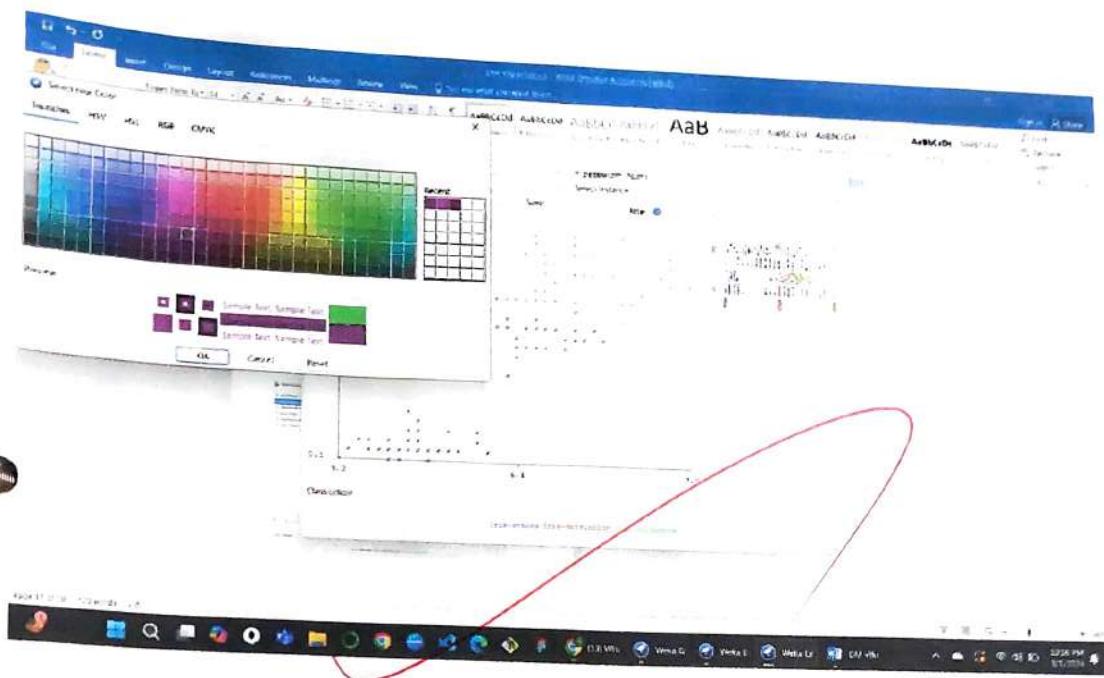


Visualizing dataset in different scenarios: **Dataset visualization** • x: sepal length & y: petal width.

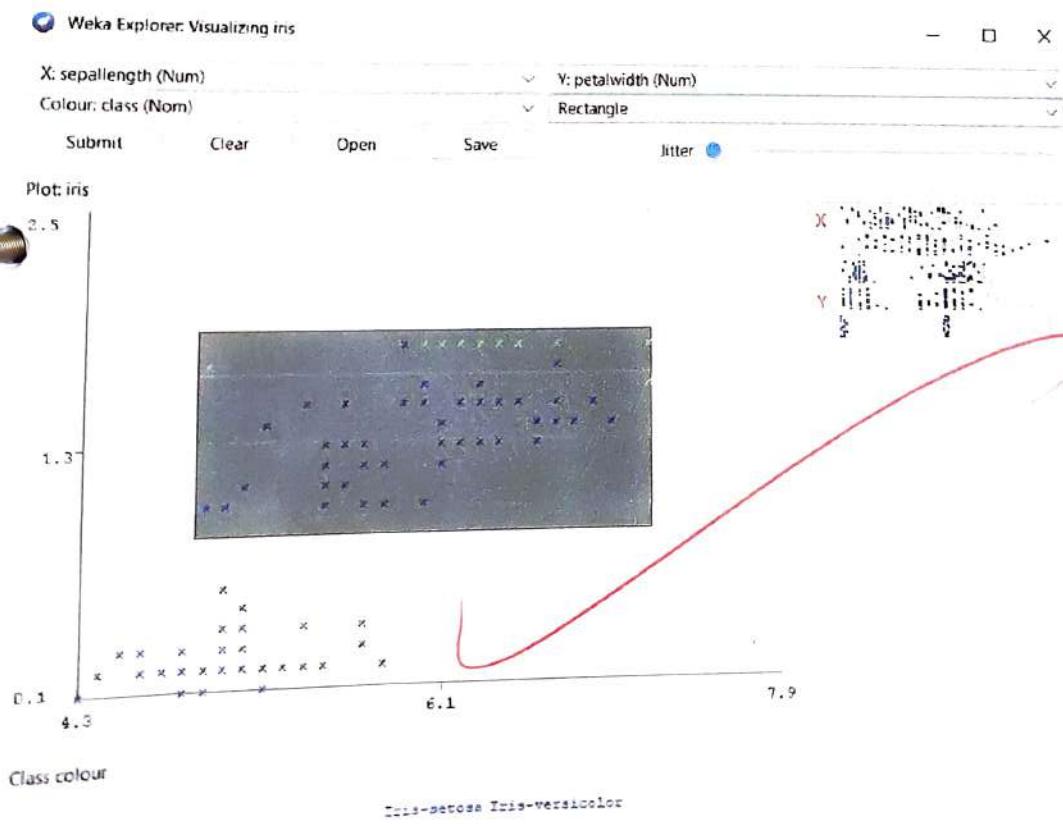


Exploration of colour tool:

18

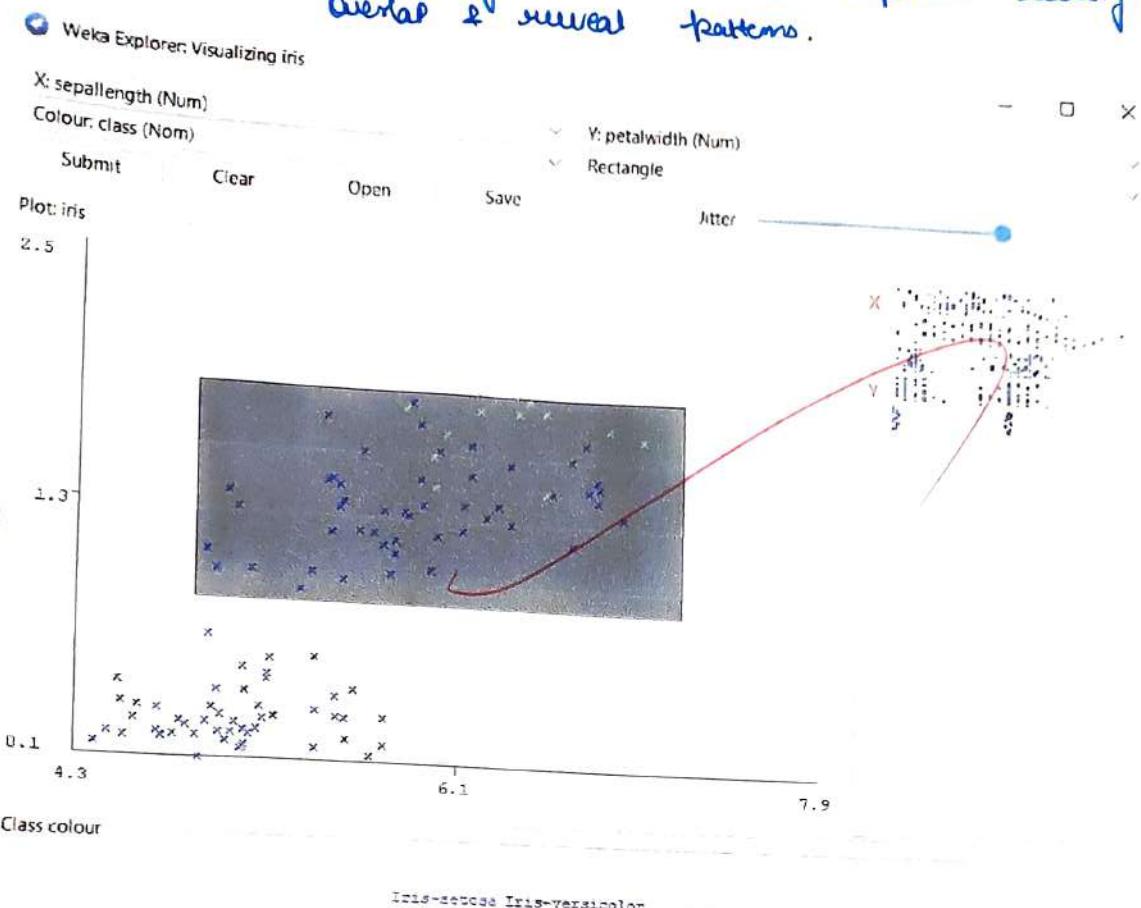


Selecting Instances: select "Rectangle" Option.



Changing the jitter:

increasing jitter, we can improve clarity, Reduce overlap & reveal patterns.



Learning Outcome:

1. understood the concept of Data visualization in Data mining.
2. learnt the process of data visualization in WEKA.

~~Vaksil Patel 20/09/24~~

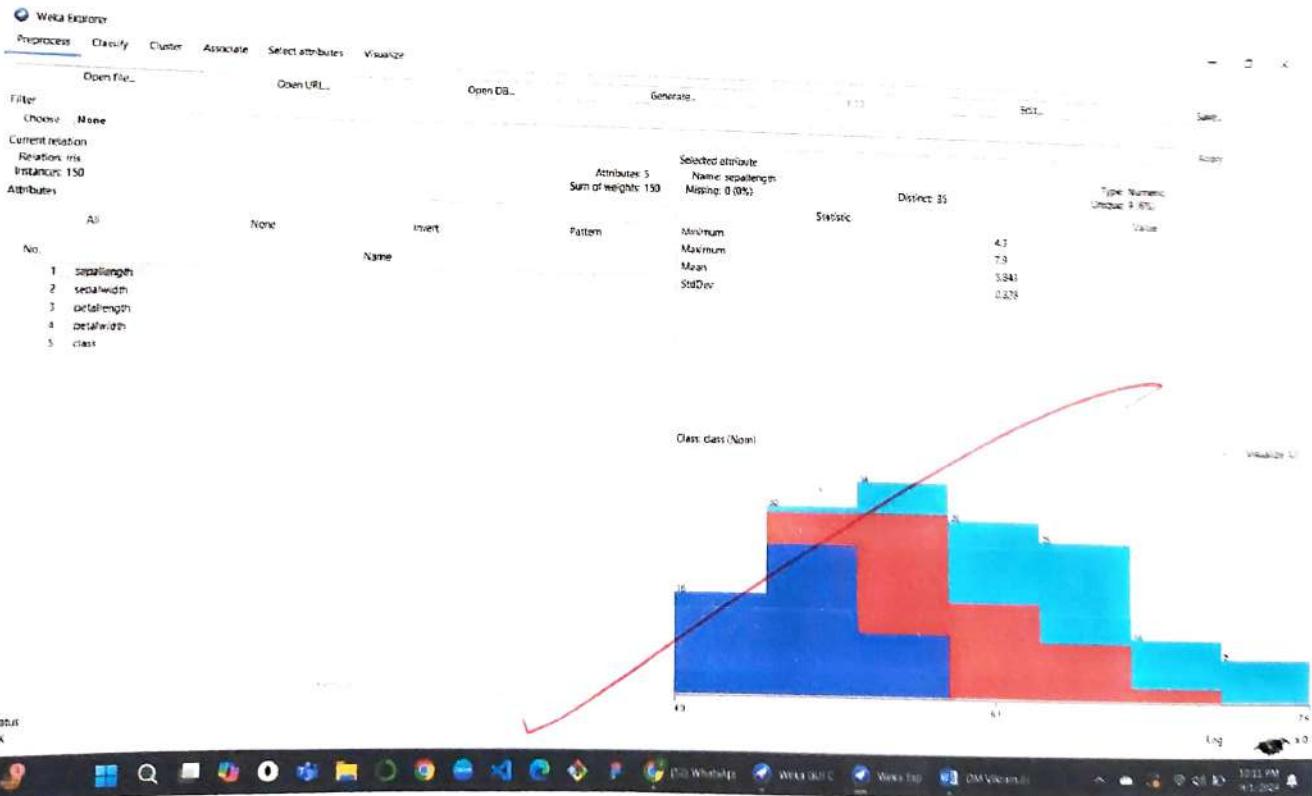
Experiment 5

Aim: Apply Clustering on a dataset using Weka.

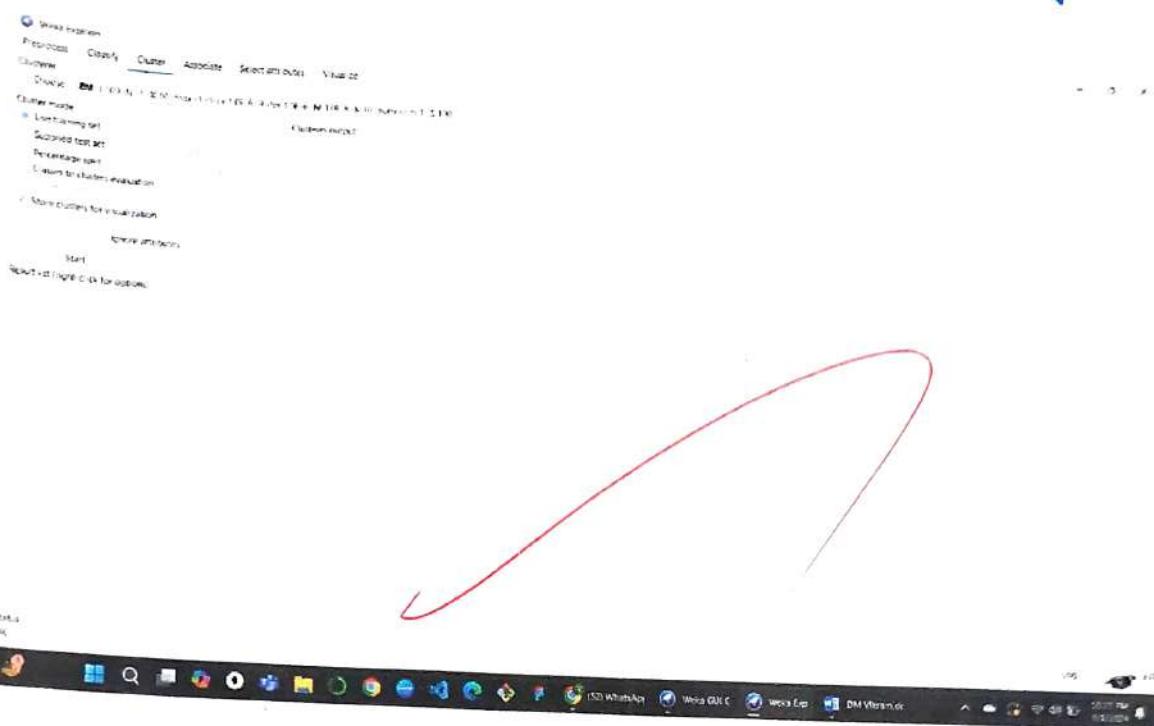
Theory: Clustering is a technique in data analysis and machine learning that groups a set of objects into clusters based on their similarities. Objects within the same cluster are more similar to each other than to those in other clusters.

- Algorithms for clustering :-
- K-means
- Hierarchical clustering
- DBSCAN
- Purpose
- Pattern Discovery
- Data Reduction.

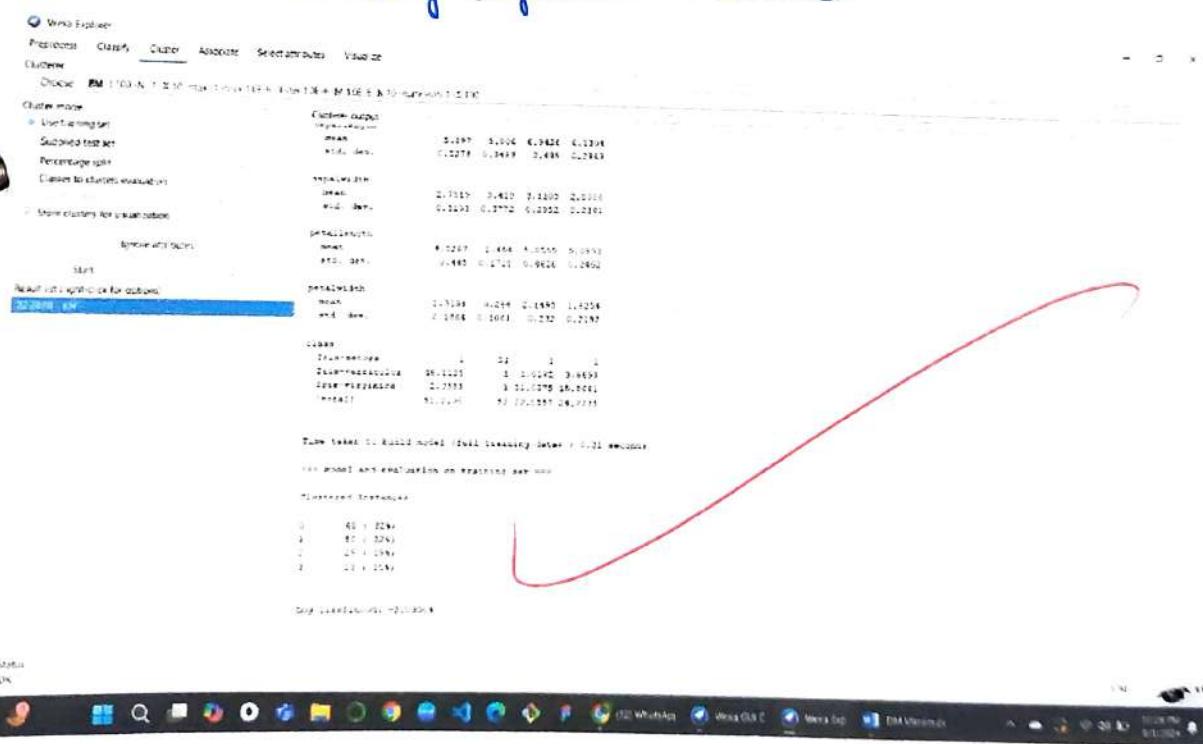
Open Iris Dataset:



Selecting cluster option: at the top of explorer dashboard of weka.



Applied clustering: from 'choose' selected clustering option, & 'start' button to apply algorithm to dataset.



Changed numCluster: Parameter of clustering algorithm, refers to no. of clusters.

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Cluster EM - 100 -N 5 -K 10 -max-1 -l 0 -cv 1 -SE -6 -inter 1 -DE -6 -M 1 -OE -6 -A 10 -min -size 1 -S 100

Cluster mode:

- Use training set
- Subsample test set
- Percentage test
- Crosses to cluster evaluation

Store clusters for visualization

Ignore attributes

Start

Result list right-click for options:
222959 - EM
222959 - EM

Clusterer output

Attribute	Class
petalLength	0.2
petalWidth	0.1
sepalLength	0.1
sepalWidth	0.1
class	0.1
iris-versicolor	0.1
iris-virginica	0.1
(total)	0.1

Simple EM (expectation maximization) class

Attributes

Class

debug: False

displayModelOldFormat: False

useNaNCheckCapability: False

iterations: 100

minimumNumberOfClusters: 1

minLogLikelihoodImprovementCV: 1.0E-6

minLogLikelihoodImprovementIterating: 1.0E-6

minStdDev: 1.0E-6

numClusters: 3

numIterationSlots: 1

numTol: 10

numXMeansRun: 10

seed: 100

Open... Save... OK Cancel

TIME taken to build model: 0m 0s

www Model and evaluation

Clustered Instances

	0	1	2
0	30 (33%)	24 (26%)	22 (23%)
1	24 (26%)	30 (33%)	22 (23%)
2	22 (23%)	22 (23%)	30 (33%)

Status: OK

Output after changing numcluster:

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Cluster EM - 100 -N 5 -K 10 -max-1 -l 0 -cv 1 -SE -6 -inter 1 -DE -6 -M 1 -OE -6 -A 10 -min -size 1 -S 100

Cluster mode:

- Use training set
- Subsample test set
- Percentage test
- Crosses to cluster evaluation

Store clusters for visualization

Ignore attributes

Start

Result list right-click for options:
222959 - EM
222959 - EM

Clusterer output

Attribute	Class
petalLength	0.2
petalWidth	0.1
sepalLength	0.1
sepalWidth	0.1
class	0.1
iris-setosa	0.1
iris-versicolor	0.1
iris-virginica	0.1
(total)	0.1

Simple EM (expectation maximization) class

Attributes

Class

debug: False

displayModelOldFormat: False

useNaNCheckCapability: False

iterations: 100

minimumNumberOfClusters: 1

minLogLikelihoodImprovementCV: 1.0E-6

minLogLikelihoodImprovementIterating: 1.0E-6

minStdDev: 1.0E-6

numClusters: 2

numIterationSlots: 1

numTol: 10

numXMeansRun: 10

seed: 100

Open... Save... OK Cancel

TIME taken to build model (full training data): 0m 0s

www Model and evaluation

Clustered Instances

	0	1
0	30 (33%)	24 (26%)
1	24 (26%)	30 (33%)
2	22 (23%)	22 (23%)
3	22 (23%)	22 (23%)

long: 0.00011840000000000001

Status: OK

Applied Hierarchical Cluster: Clustering algorithm which creates a hierarchy of clusters, either by merging smaller clusters or dividing larger ones.

```

Weka Explorer
Preprocess Classify Cluster Associate Select attributes Visualize
Cluster Choose HierarchicalCluster - N=2-L=SINGLE - A=euclideanDistances-R=firstSet
Cluster mode
  Use training set
  Supplied test set
  Percentage split
  Classes to clusters evaluation
  Store clusters for visualization
  Ignore attributes
Start
Results for right-click options:
22.28.08 - PM
22.29.10 - EM
22.30.10 - EM
22.31.12 - HierarchicalCluster
22.32.12 - HierarchicalCluster

HierarchicalCluster - N=2-L=SINGLE - A=euclideanDistances-R=firstSet
Cluster output
  repalength
  repalwidth
  getalength
  getalwidth
  class
Test mode: split 80% train, remainder test
== Clustering model (full training set) ==
Clusters:
Cluster 0: {((0.0, 0.03254, 0.9, 0.01284), (0.0, 0.03254, 0.0, 0.03254), (0.0, 0.03254, 0.0, 0.03254))} | 0.00332, ((0.1, 0.02779, 0.0, 0.02779), (0.1, 0.02779, 0.1, 0.02779), (0.1, 0.02779, 0.1, 0.02779))
Cluster 1: {((0.0, 0.03254, 0.0, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284))} | 0.00332
Time taken to build model (full training data): 0.01 seconds
== Model and evaluation on training set ==
Clustered Instances:
  0: 52 + 3384
  1: 19 + 4794

```

Output of Hierarchical Cluster:

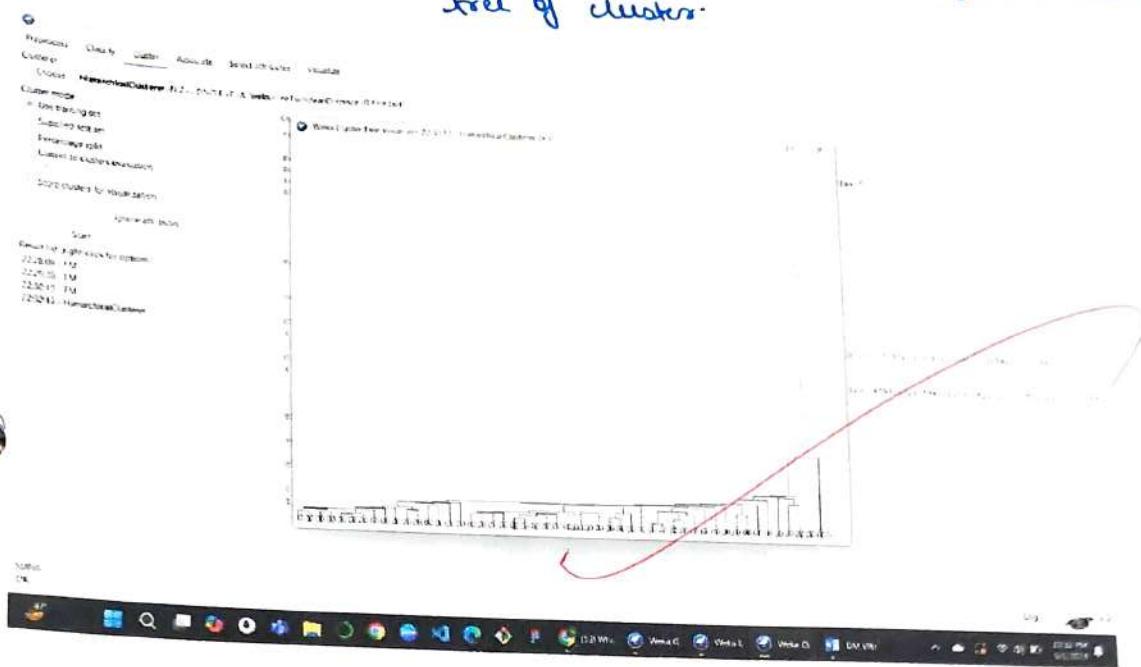
```

Weka Explorer
Preprocess Classify Cluster Associate Select attributes Visualize
Cluster Choose HierarchicalCluster - N=2-L=SINGLE - A=euclideanDistances-R=firstSet
Cluster mode
  Use training set
  Supplied test set
  Percentage split %: 80
  Classes to clusters evaluation
  Store clusters for visualization
  Ignore attributes
Start
Results for right-click options:
22.28.08 - PM
22.29.10 - EM
22.30.10 - EM
22.31.12 - HierarchicalCluster
22.32.12 - HierarchicalCluster

HierarchicalCluster - N=2-L=SINGLE - A=euclideanDistances-R=firstSet
Cluster output
  repalength
  repalwidth
  getalength
  getalwidth
  class
Test mode: split 80% train, remainder test
== Clustering model (full training set) ==
Clusters:
Cluster 0: {((0.0, 0.03254, 0.9, 0.01284), (0.0, 0.03254, 0.0, 0.03254), (0.0, 0.03254, 0.0, 0.03254))} | 0.00332, ((0.1, 0.02779, 0.0, 0.02779), (0.1, 0.02779, 0.1, 0.02779), (0.1, 0.02779, 0.1, 0.02779))
Cluster 1: {((0.0, 0.03254, 0.0, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284))} | 0.00332
Time taken to build model (full training data): 0.02 seconds
== Model and evaluation on test split ==
Cluster 0: {((0.0, 0.03254, 0.9, 0.01284), (0.0, 0.03254, 0.0, 0.03254), (0.0, 0.03254, 0.0, 0.03254))} | 0.00332, ((0.1, 0.02779, 0.0, 0.02779), (0.1, 0.02779, 0.1, 0.02779), (0.1, 0.02779, 0.1, 0.02779))
Cluster 1: {((0.0, 0.03254, 0.0, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284), (0.0, 0.03254, 0.1, 0.01284))} | 0.00332
Time taken to build model (percentage split): 0.01 seconds
Clustered Instances:
  0: 19 + 619
  1: 11 + 2794

```

Visualization of Hierarchical Cluster Tree: choose tree visualize, to visualize tree of clusters.



Learning Outcome:

- Understood the clustering concept in data mining.
- learnt applying different clustering options on dataset in weka.
- Visualize the clustering tree. ✓

last slide 23/09/24

Experiment 6

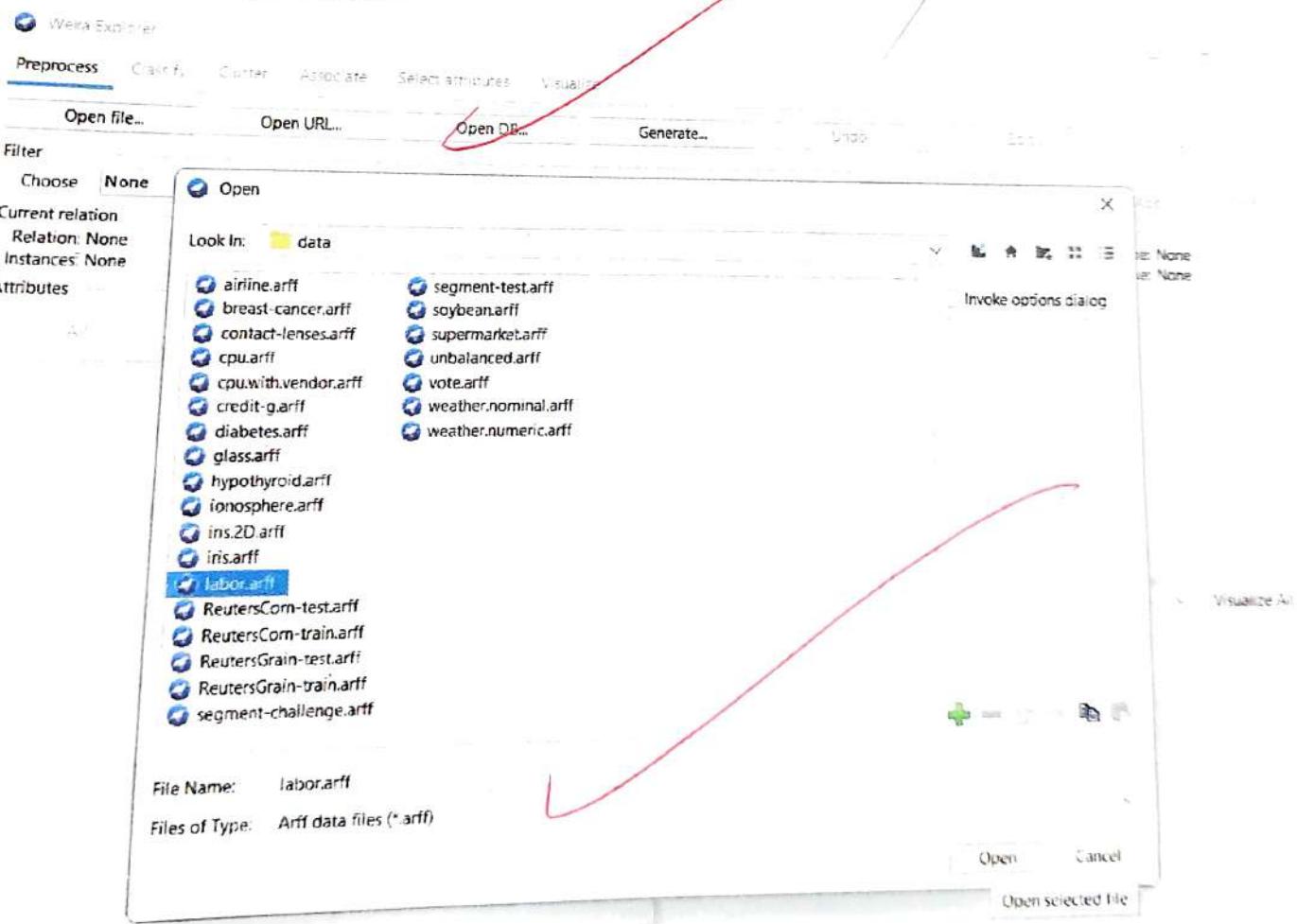
Aim: Apply Classification on a dataset using Weka.

Theory: Classification is a supervised machine learning technique used to predict the category or class of a given data point based on its features.

It involves a training model on labeled data to learn patterns and relationships, which can then be used to classify new, unseen data points.

- Common algorithms are:-
- Decision Tree
- Naive Bayes
- Support Vector Machine

Selection of Labor dataset:

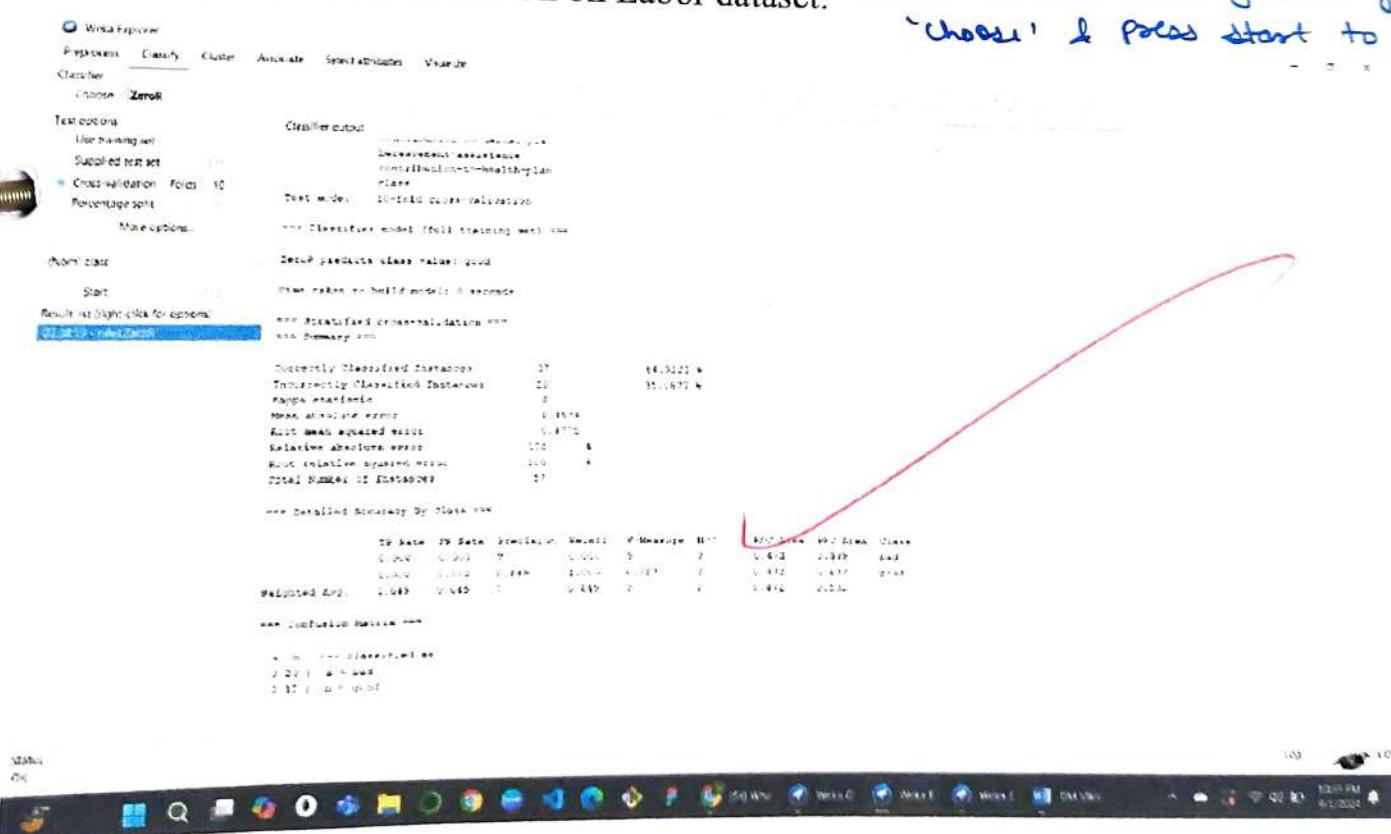


Analyzation of Labor Dataset:



Applying ZeroR Classification on Labor dataset:

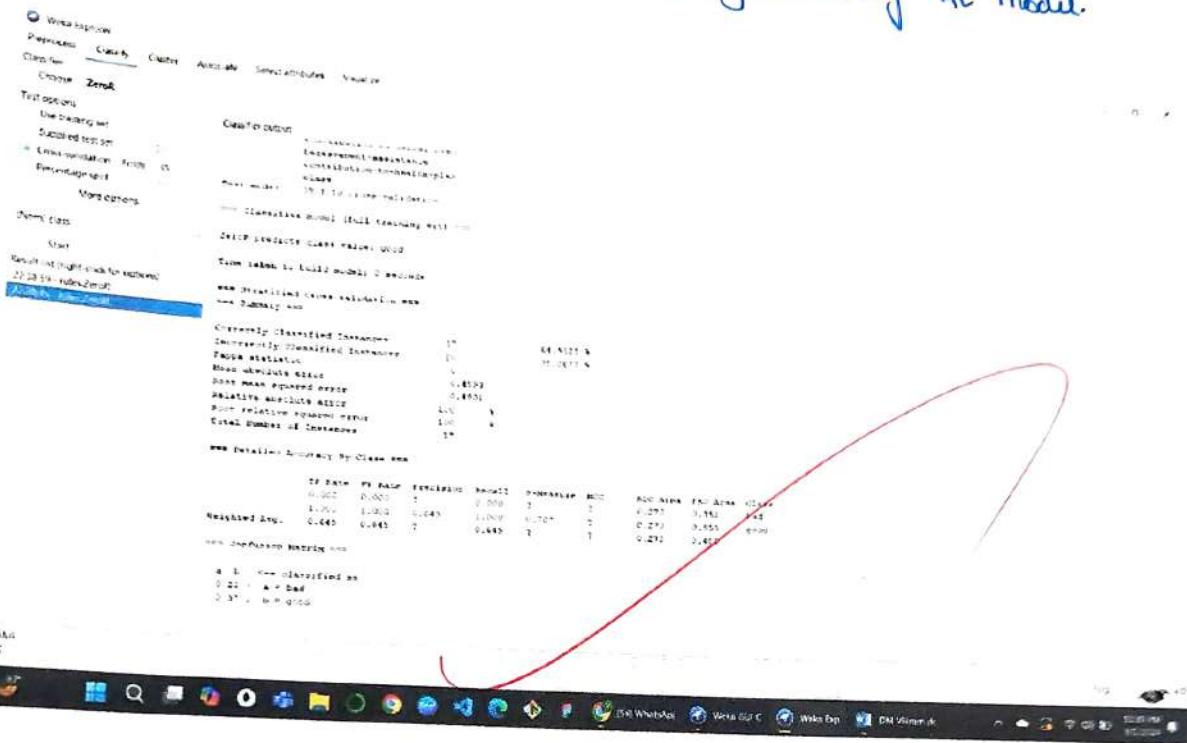
select ZeroR classification from 'choose' & press start to apply.



Changing number of Cross-Validation Folds:

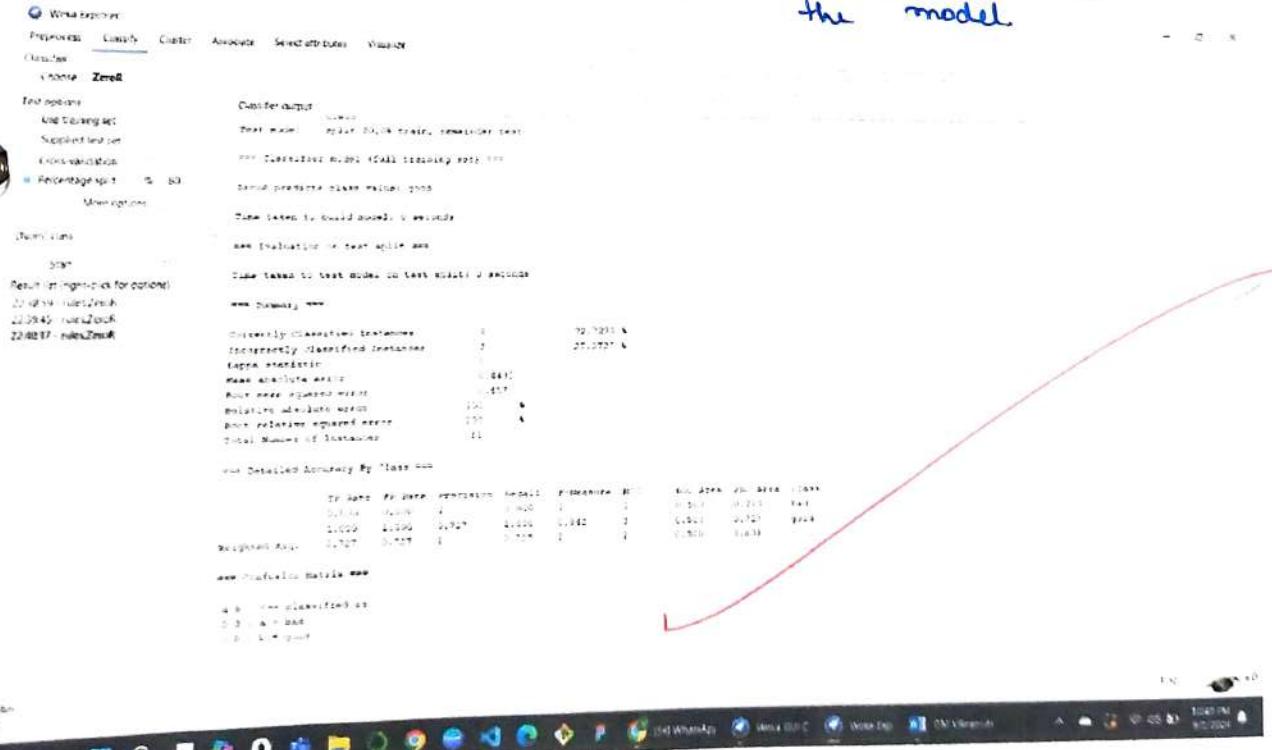
27

ss-validation Folds: refers to subsets of a dataset to evaluate the performance of ML model.



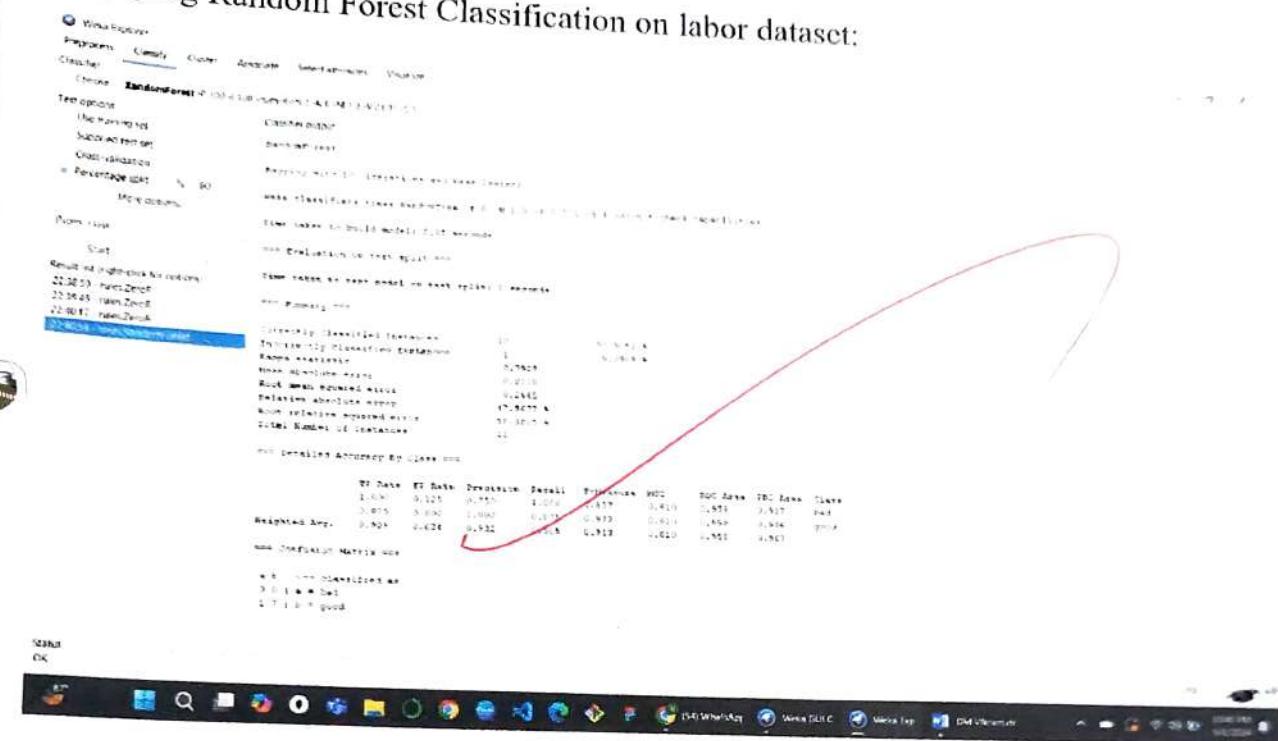
Splitting the dataset into training and testing sets:

80% for training & 20% for testing
the model.

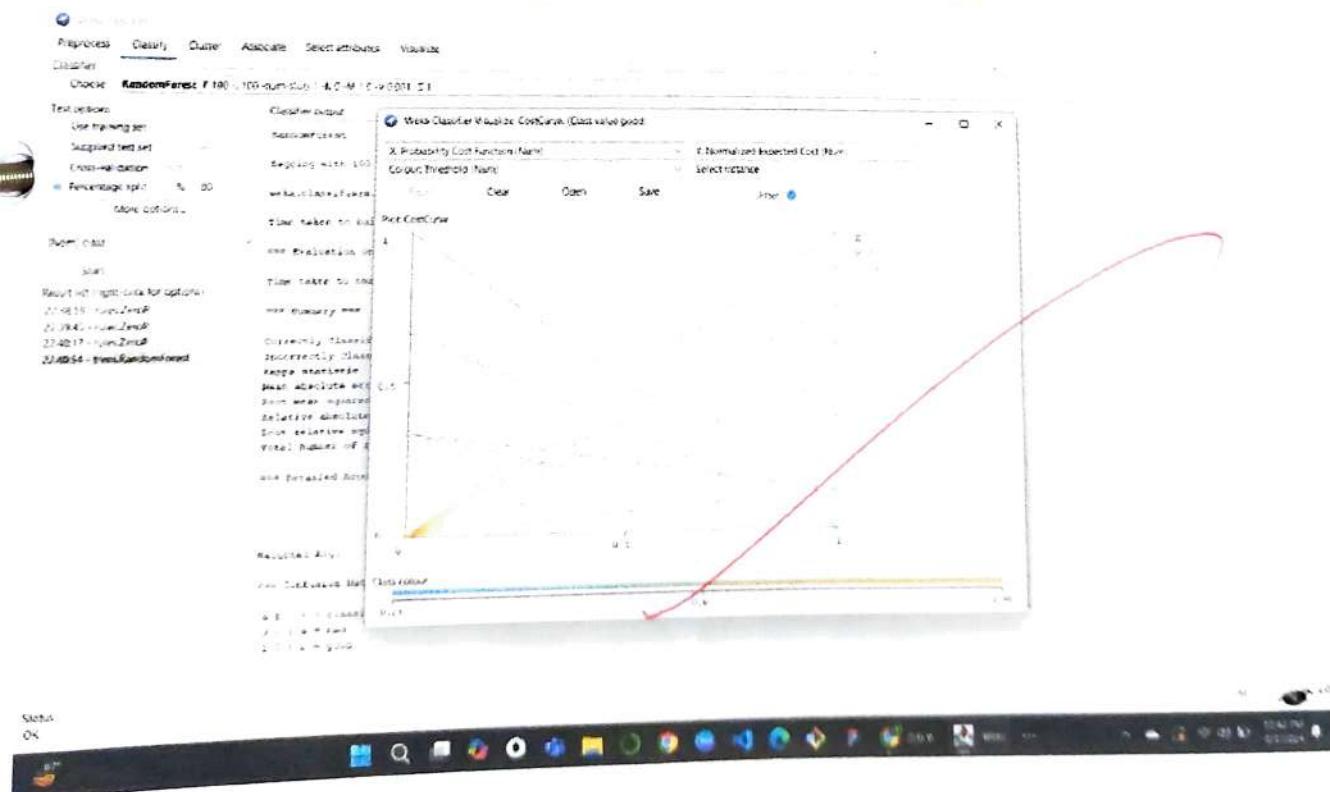


We can choose multiple classifiers in weka, Random forest is one of them, which is an ensemble learning method.

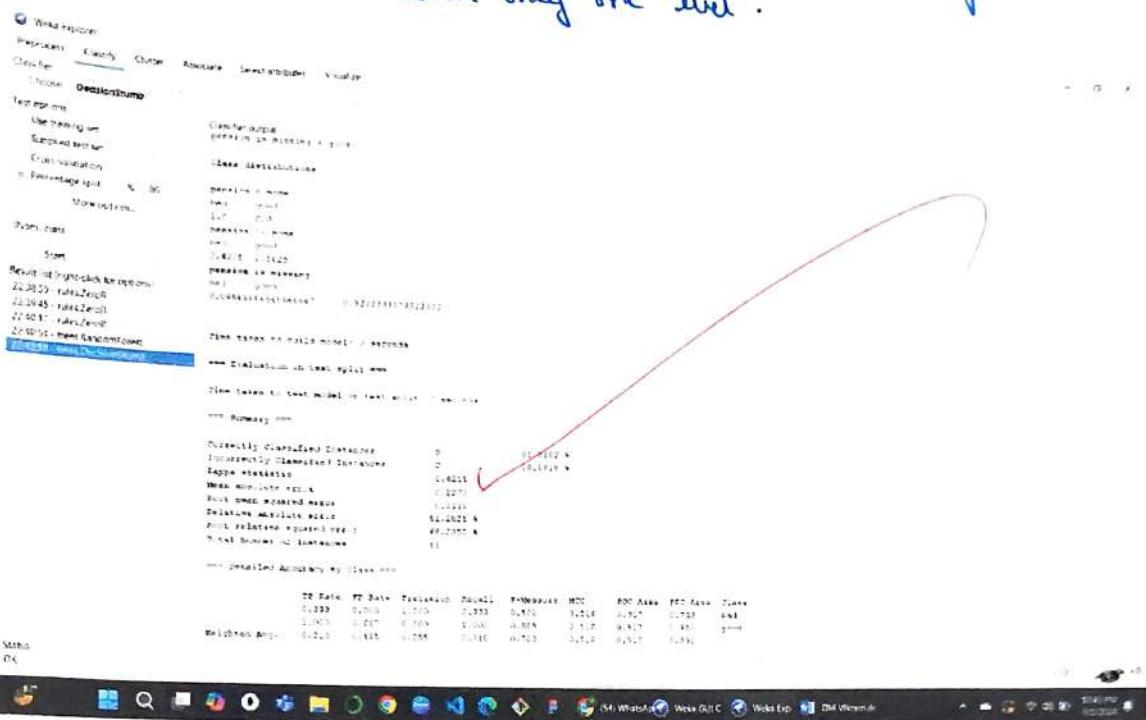
Applying Random Forest Classification on labor dataset:



Visualization of cost curve: Performance of model analysis .



Applying Decision Stump Classification on labor dataset: uses only one decision tree with only one level.



Learning Outcome:

- understood the concepts of classification.
- learnt how to classify datasets and visualize performance of model in weka.

Completed
23/09/24

Experiment 7

Objective: Evaluate the performance of classification techniques using different parameters.

Theory:

Classification is a supervised learning technique in data mining and machine learning where a model is trained on labeled data to categorize new instances into classes.

In this experiment, we aim to explore and compare different classification techniques such as decision trees, Naive Bayes, SVM and KNN.

Key classification techniques

(I) Decision Trees :-

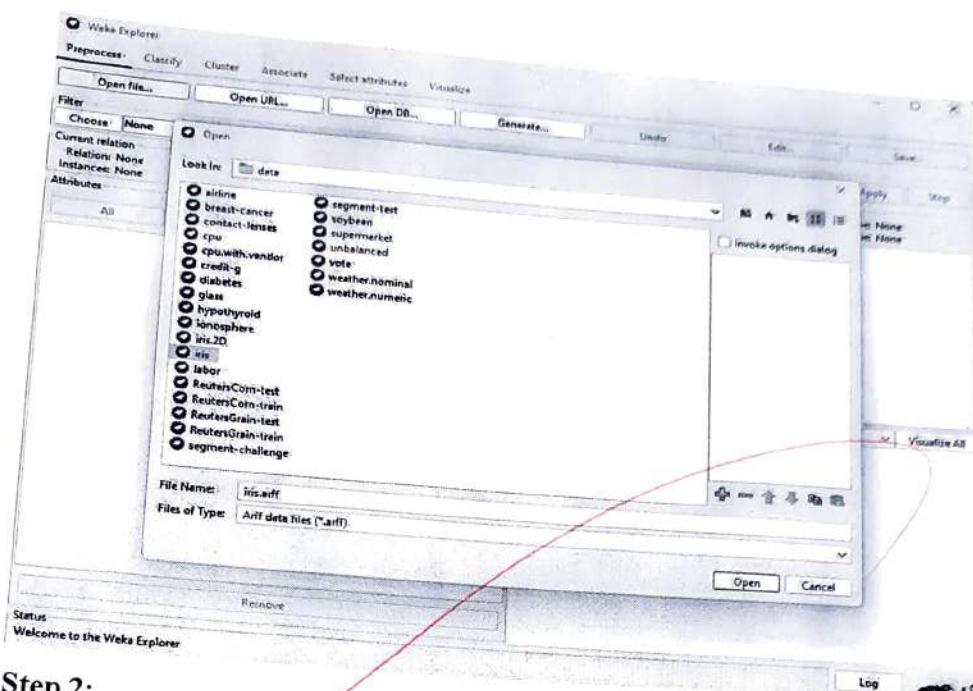
A tree structured classifier where each node represents a feature, each branch represents a decision rule and each leaf node represents an outcome.

II) Naive Bayes :- A probabilistic classifier based on Baye's theorem assuming independence between features.

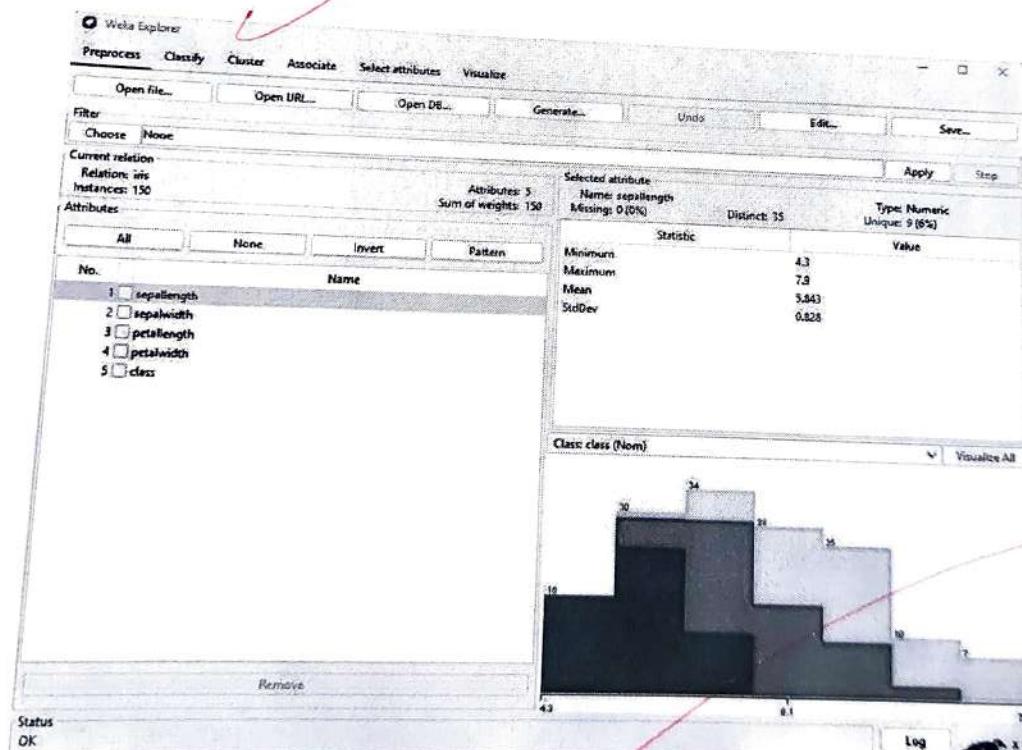
III) Support Vector Machine (SVM) :- SVM finds a hyperplane that separates classes. It is particularly powerful for binary classification.

IV) K - Nearest neighbours :- A non-parametric method where classification is based on the majority label of the nearest neighbour.

Step 1:



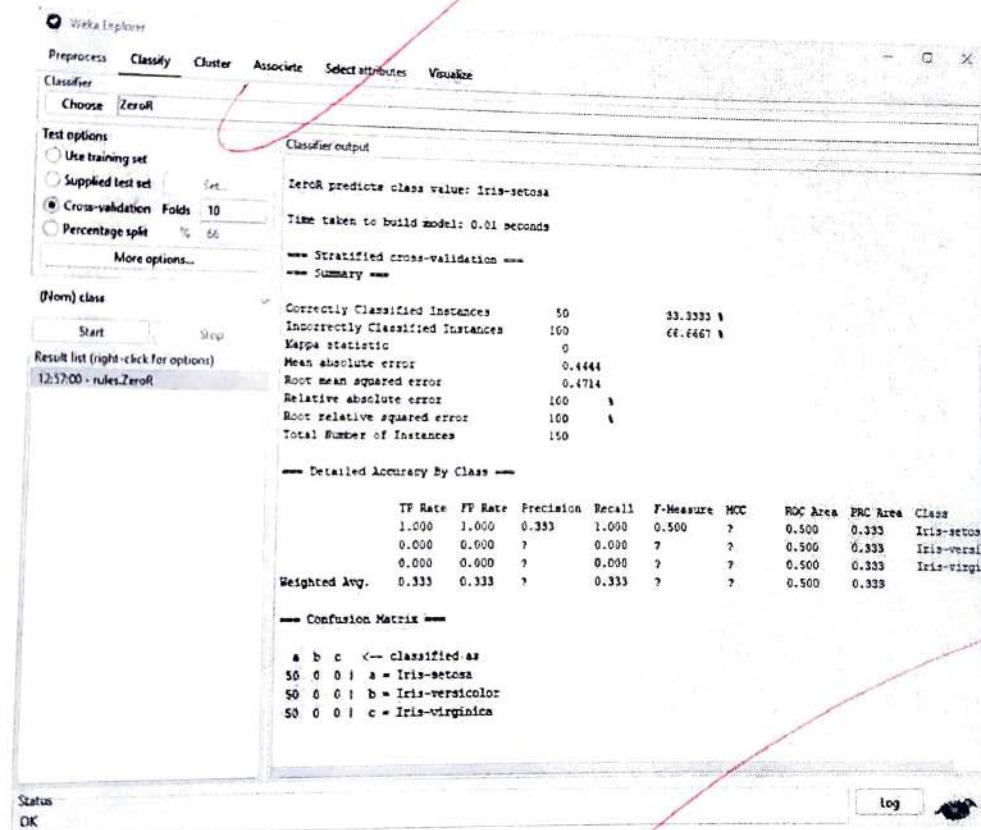
Step 2:



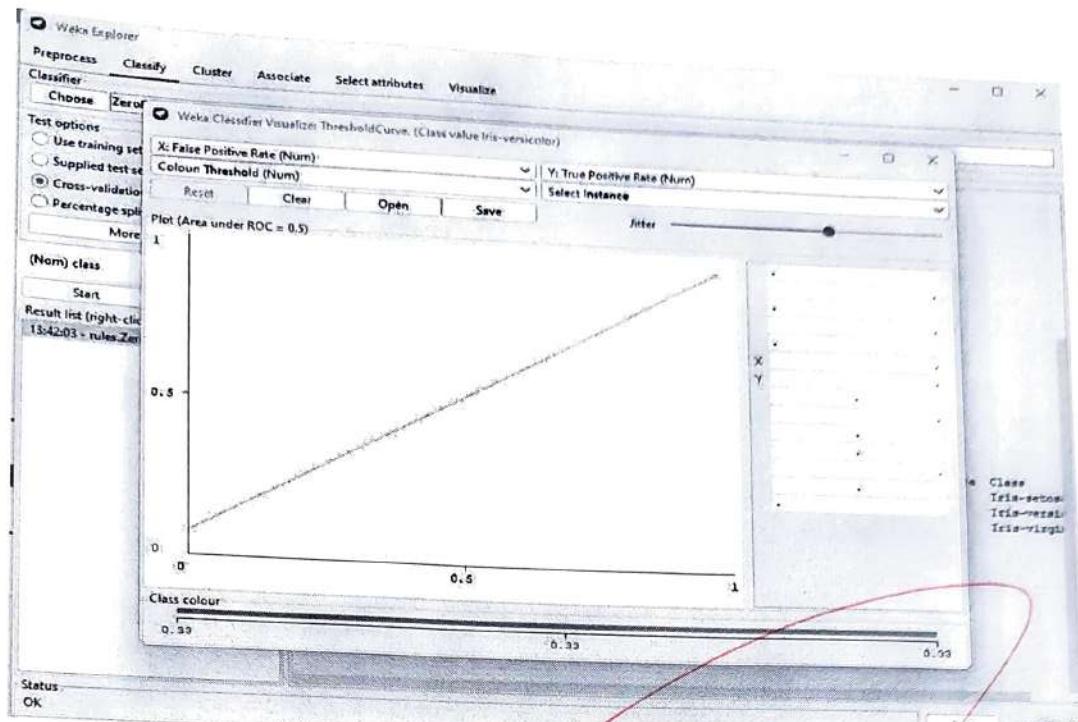
Step 3:



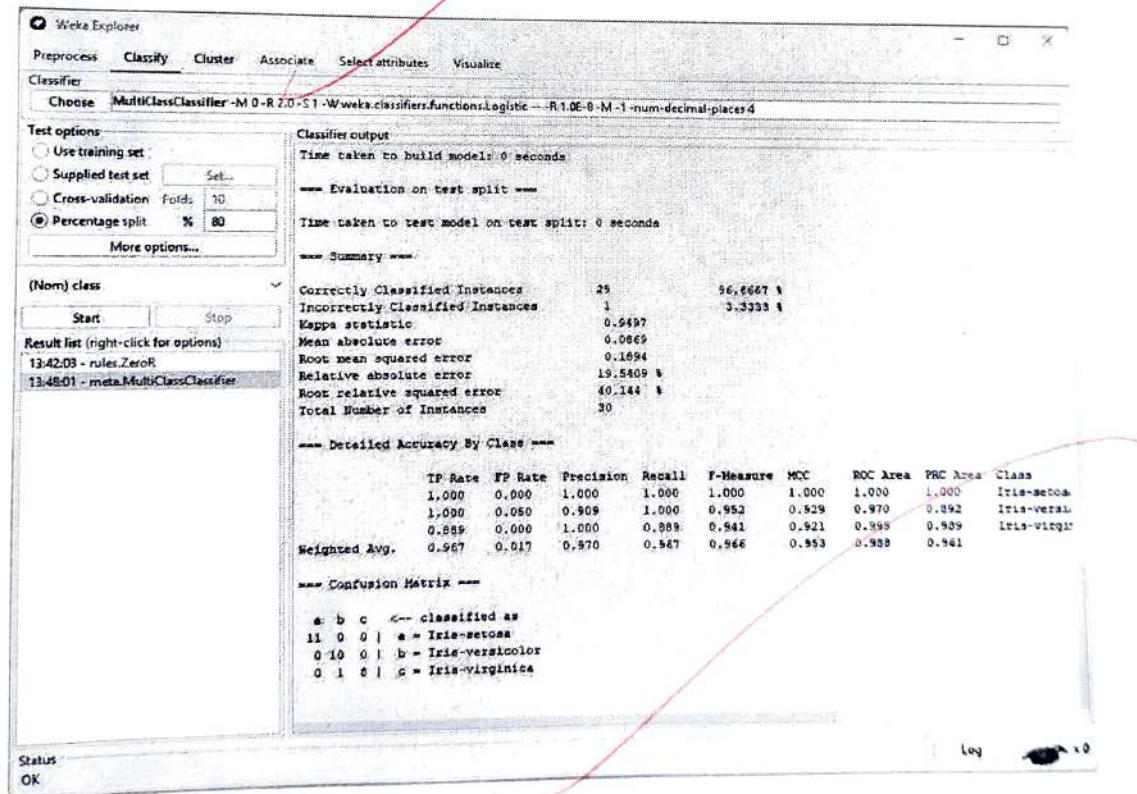
Step 4:



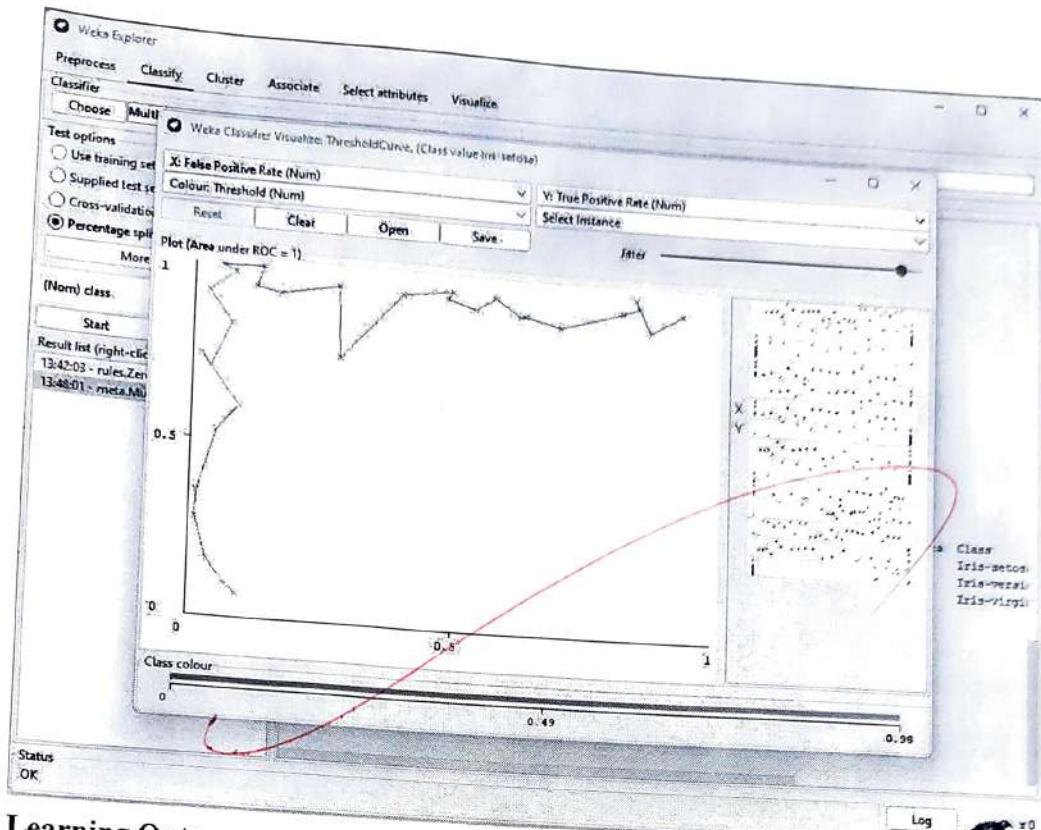
Step 5 :



Step 6:



Step 7:



Learning Outcomes:

- understood various classification techniques like SVM, decision trees, naive Bayes, KNN.
- learnt to implement & evaluate performance of classification techniques using different parameters.

Aloftch Shubh

Experiment 8

Objective: Implementation of Bagging and Boosting techniques on ARFF files using WEKA.

Theory:

Bagging :- It's an ensemble learning method that improves the stability and accuracy of the algorithms. It involves creating multiple subsets of training data by random sampling, training a separate model on each subset and then averaging the prediction for regression and majority voting for classification.

Boosting :- Boosting is another ensemble technique that aims to reduce bias and improve predictive performance by sequentially training models.

ARFF files :- [Attribute Relation file format] files stored in a format, readable by the weka.

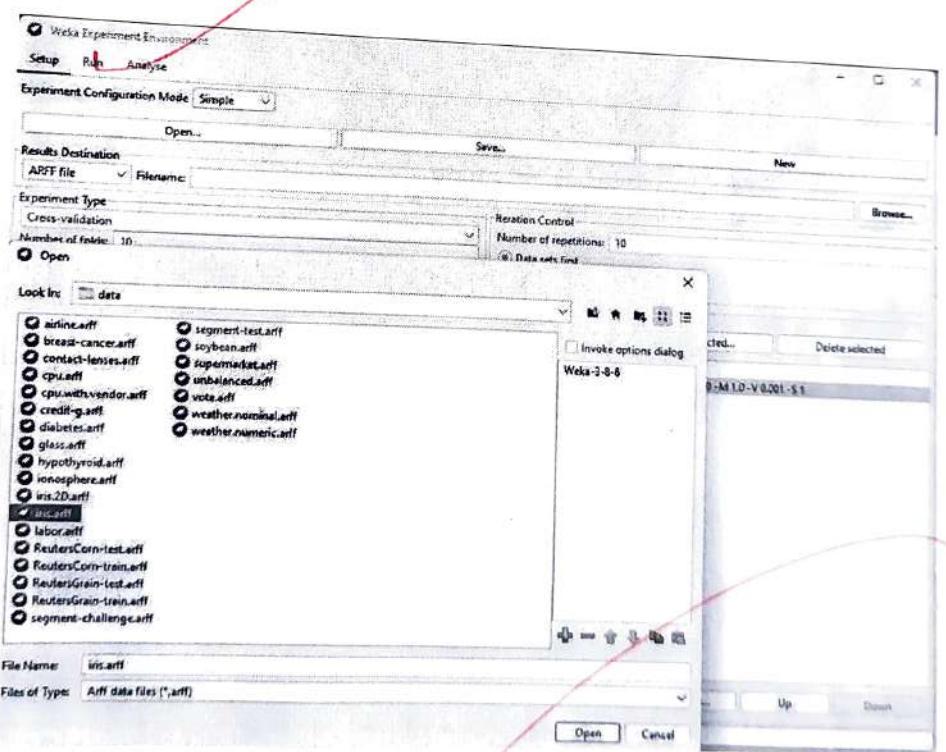
Procedure:-

- 1) Load ARFF files
- 2) Choose bagging / boosting under classify tab, select an algorithm like decision trees as a meta classifier.
- 3) Adjust parameters such as no. of iteration or sample size.
- 4) Execute the algorithm and review the performance.

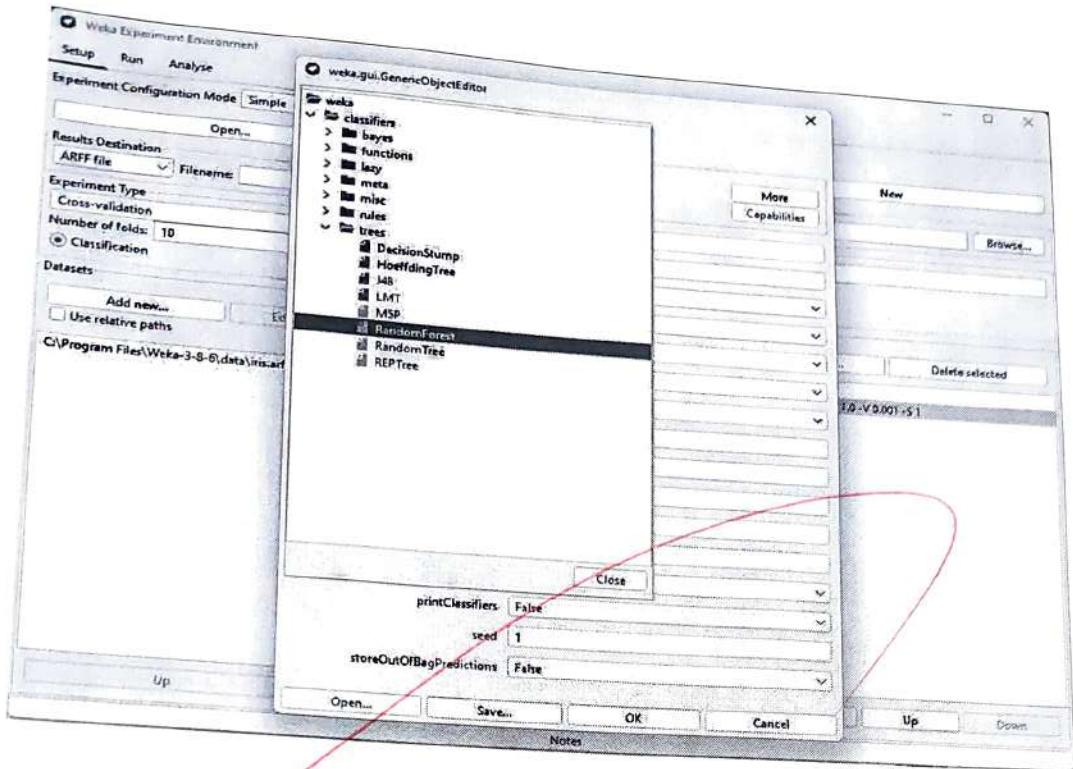
Step 1:



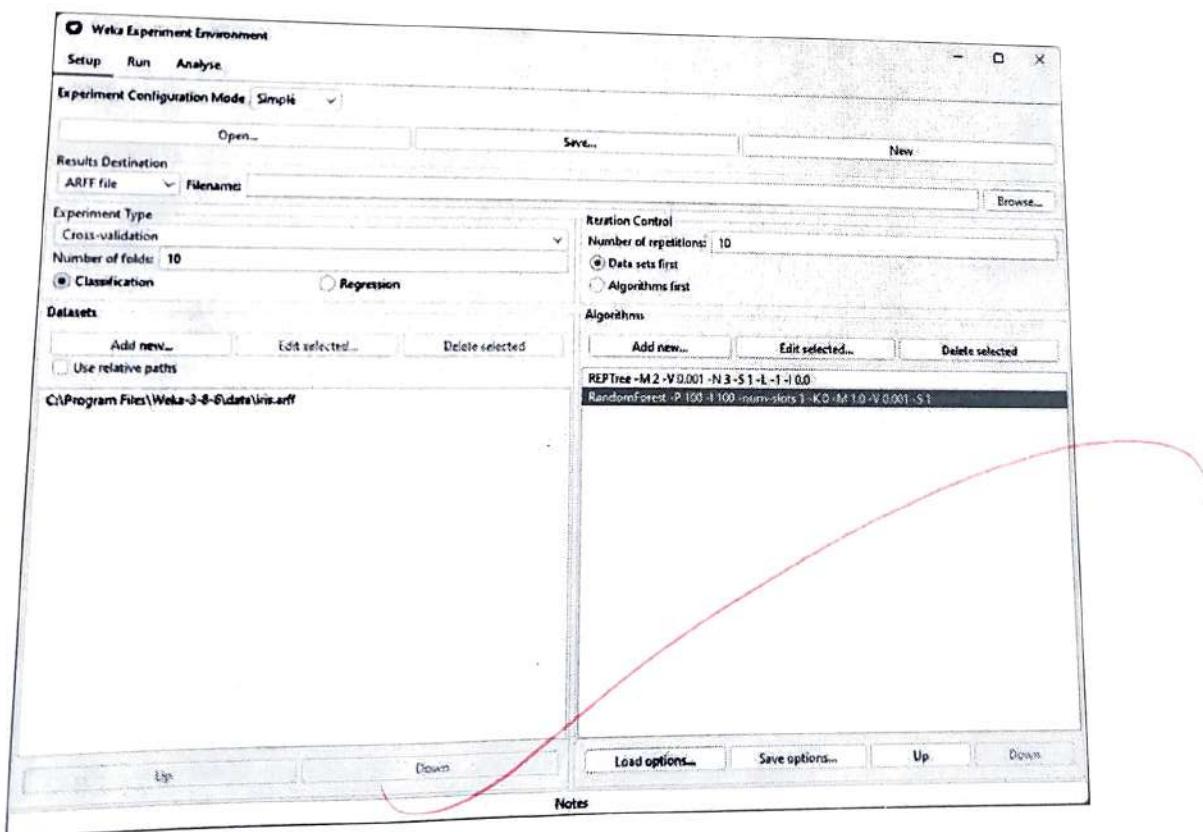
Step 2:



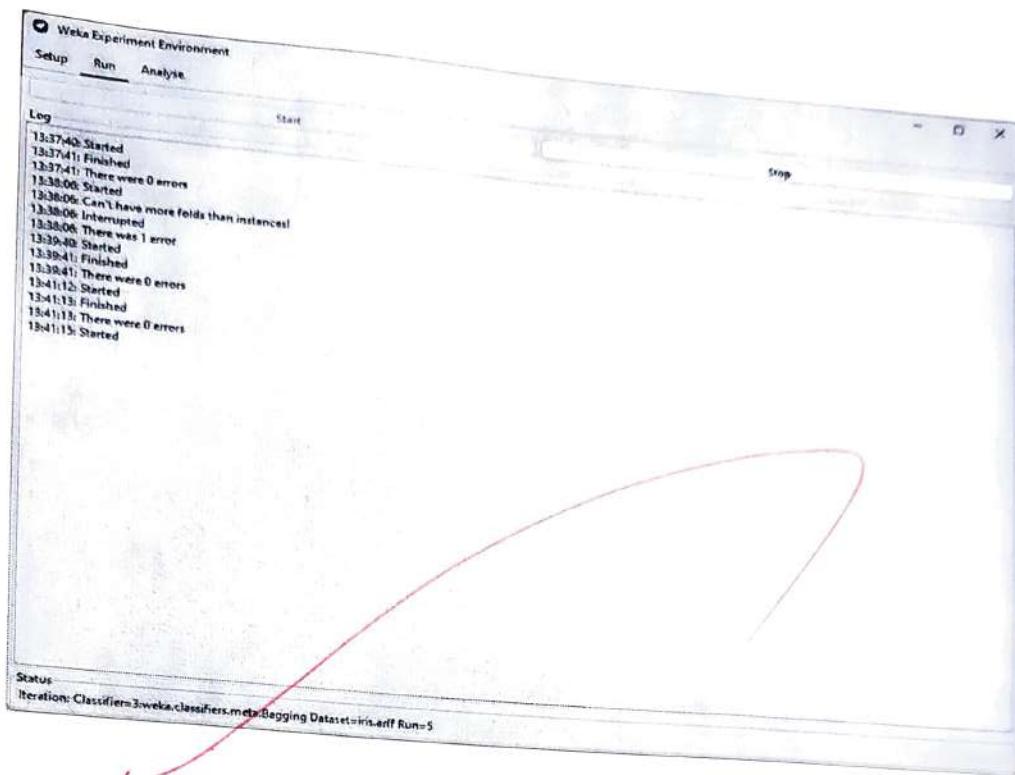
Step 3:



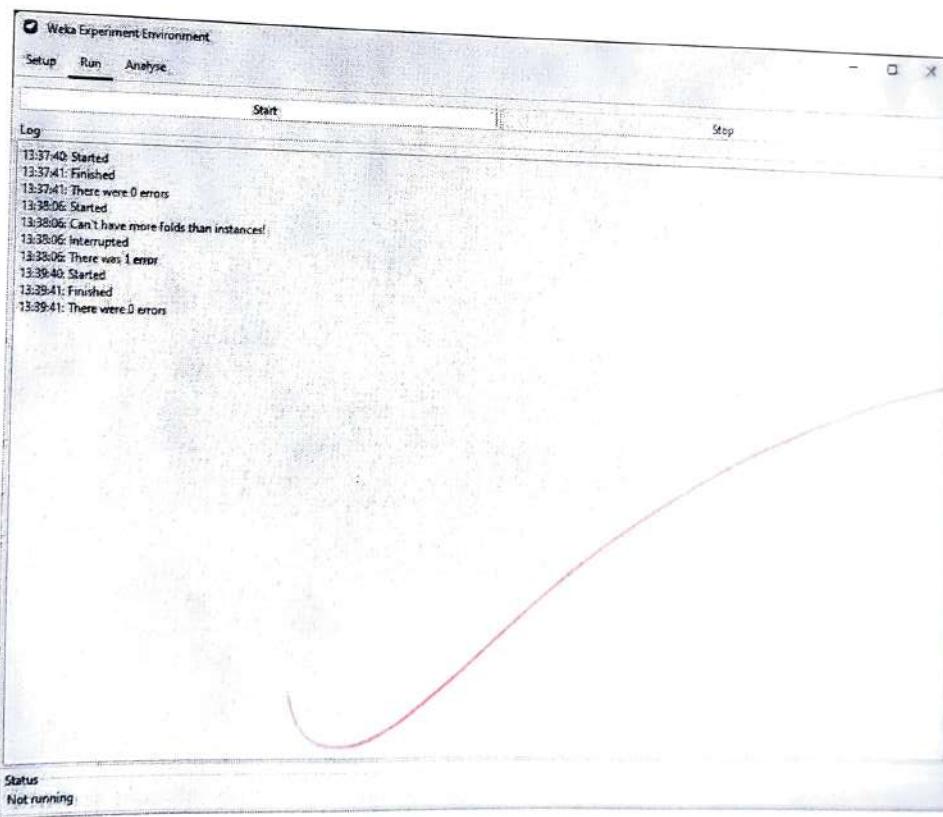
Step 4:



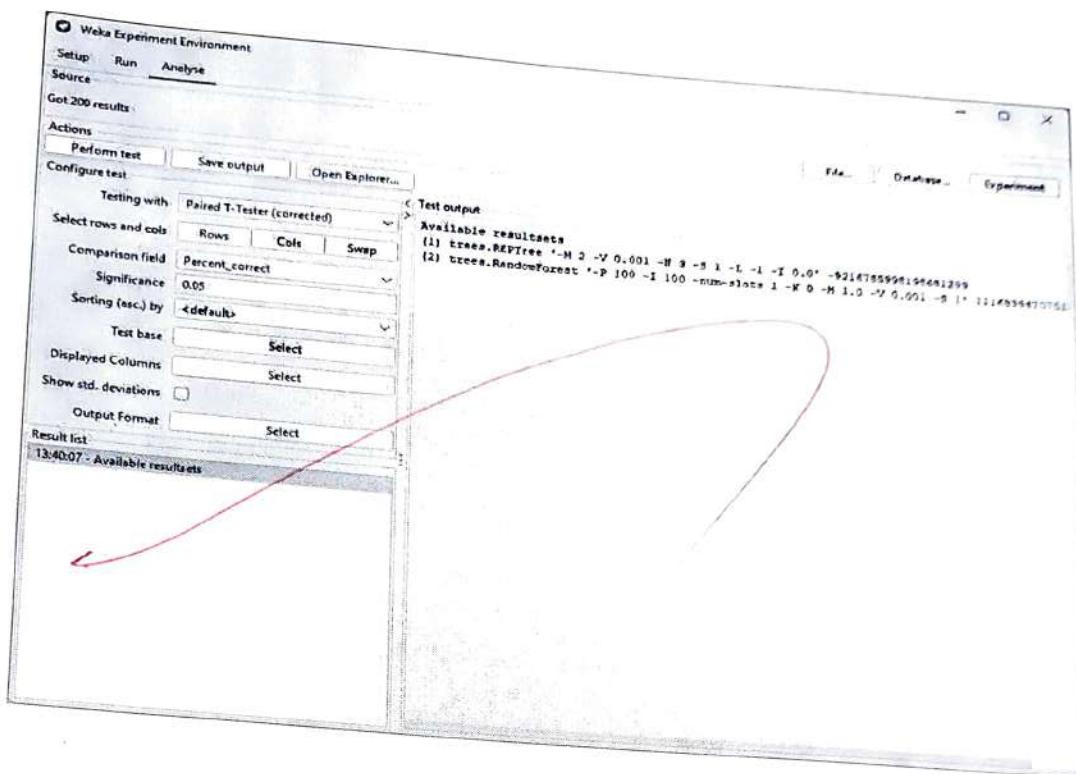
Step 5:



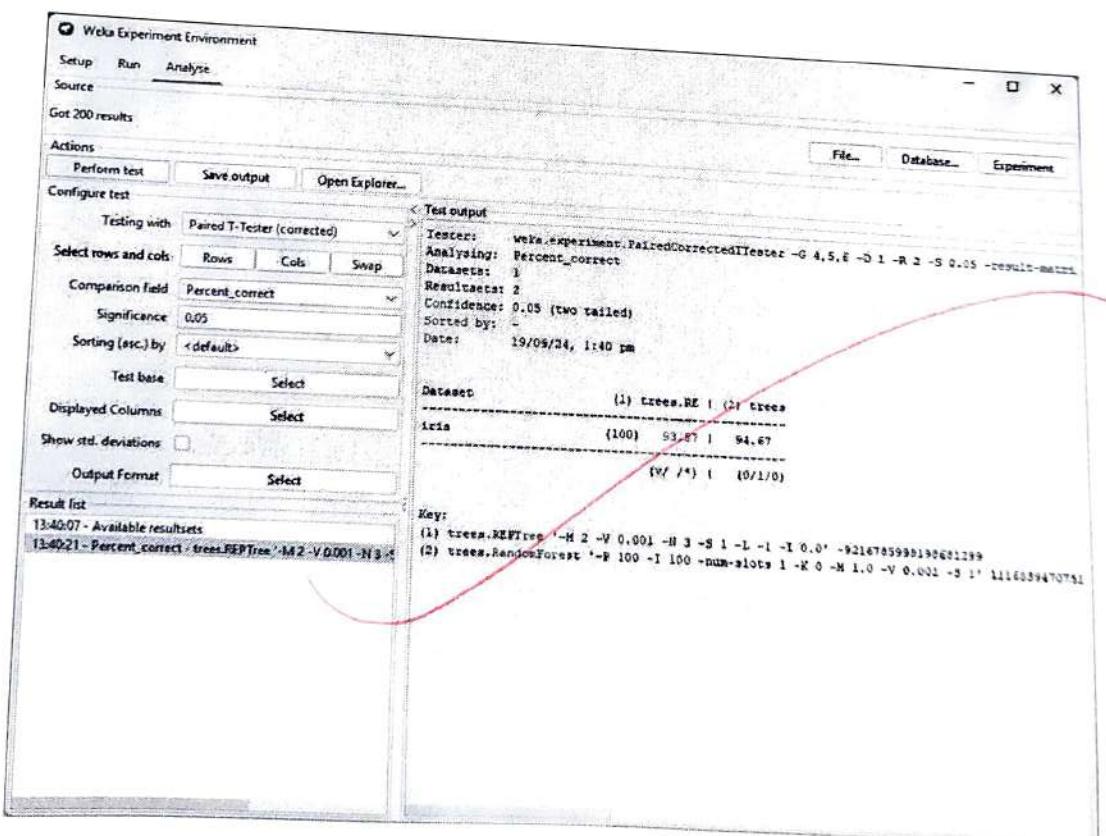
Step 6:



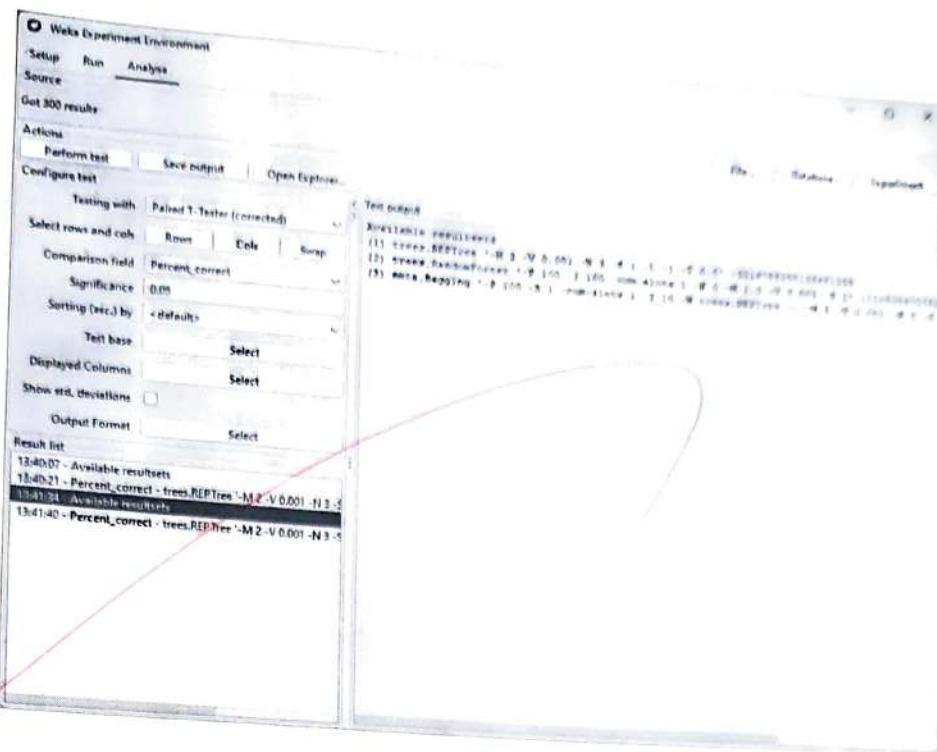
Step 7:



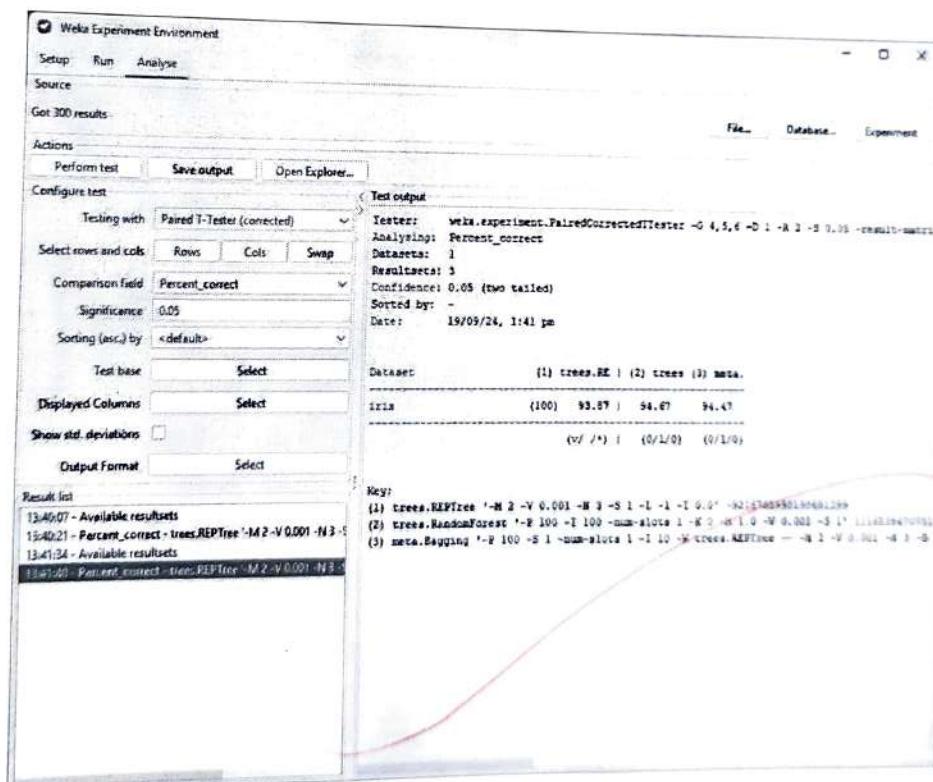
Step 8:



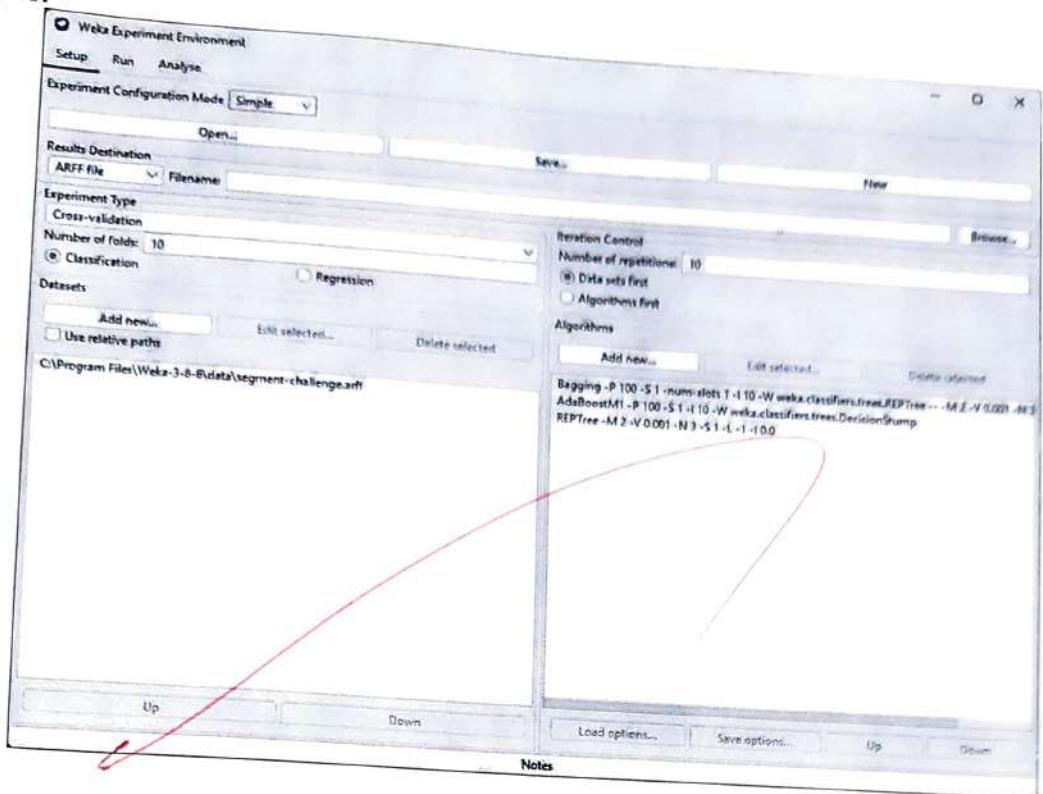
Step 9:



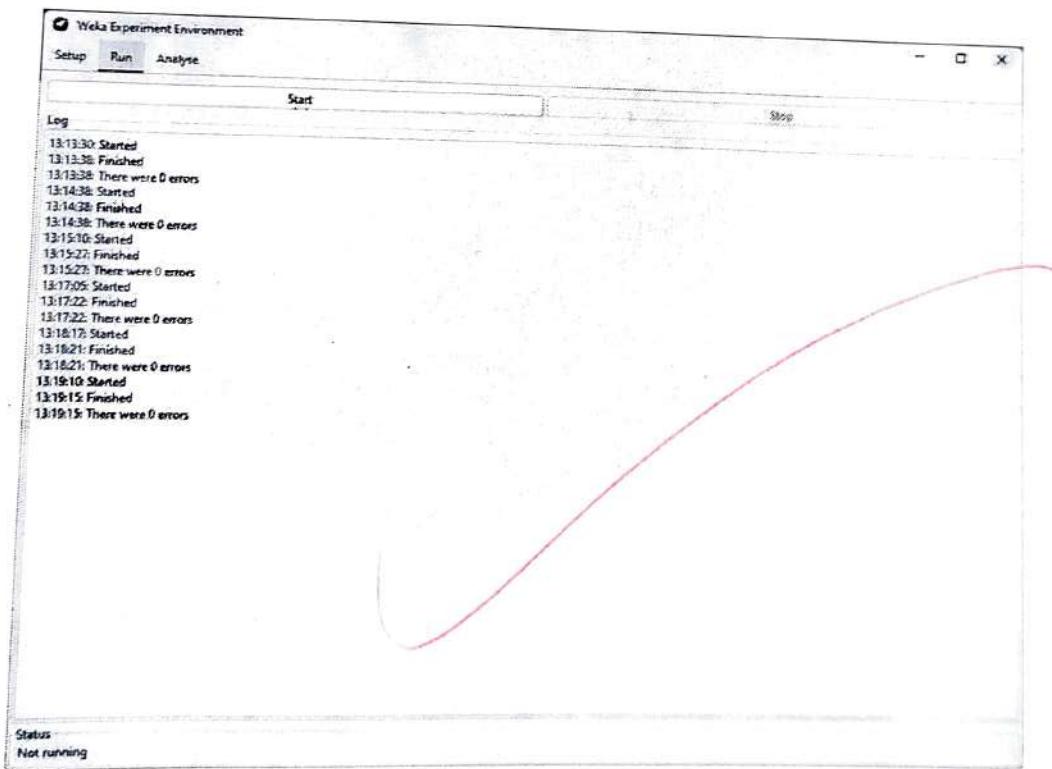
Step 10:



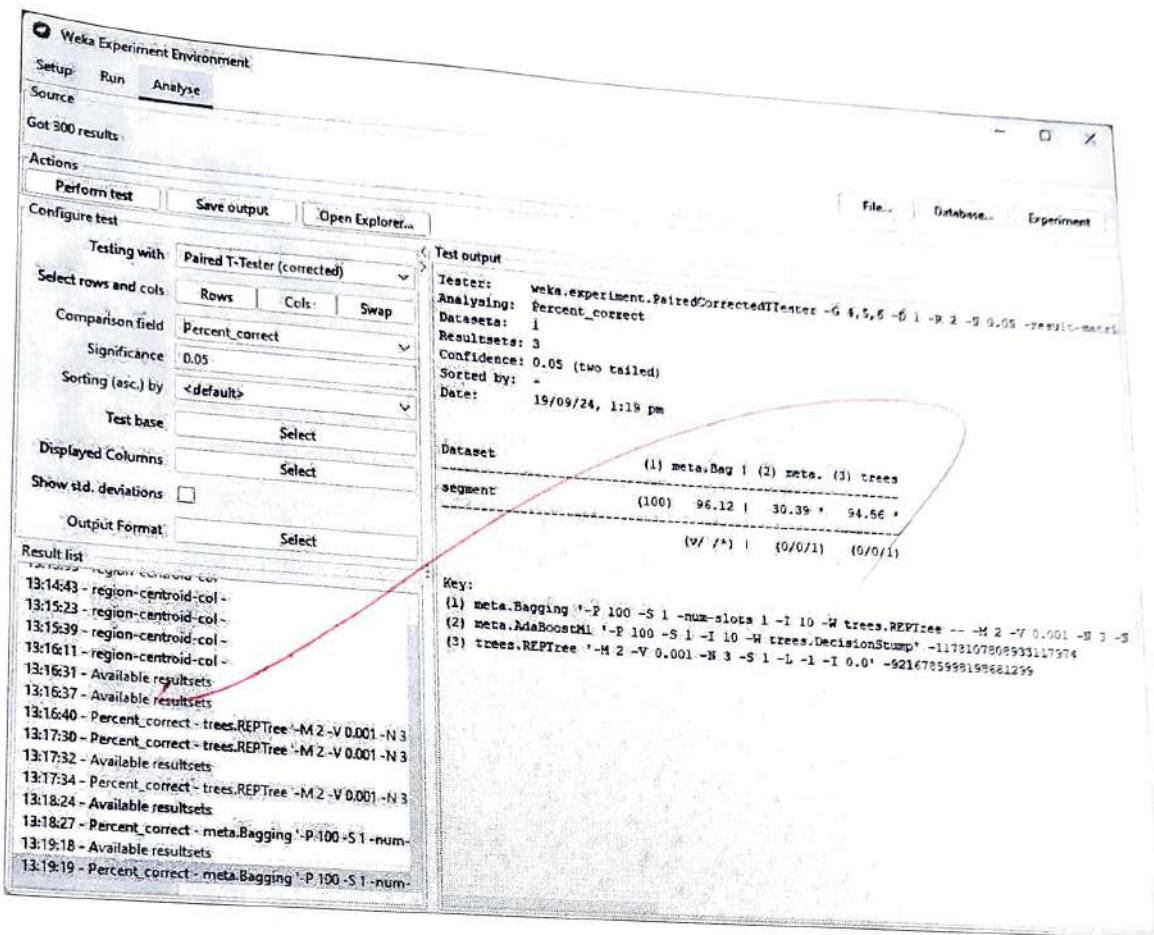
Step 11:



Step 12:



Step 13:



Learning Outcomes:

learnt how ensemble techniques can improve classification accuracy, reduce variance (bagging) and address misclassification bias (boosting).

WAPSIKUL JUL 24

Experiment 9

Objective: Apply the concept of Voting ensemble method to ARFF files and compare the results with single classifiers.

Theory:

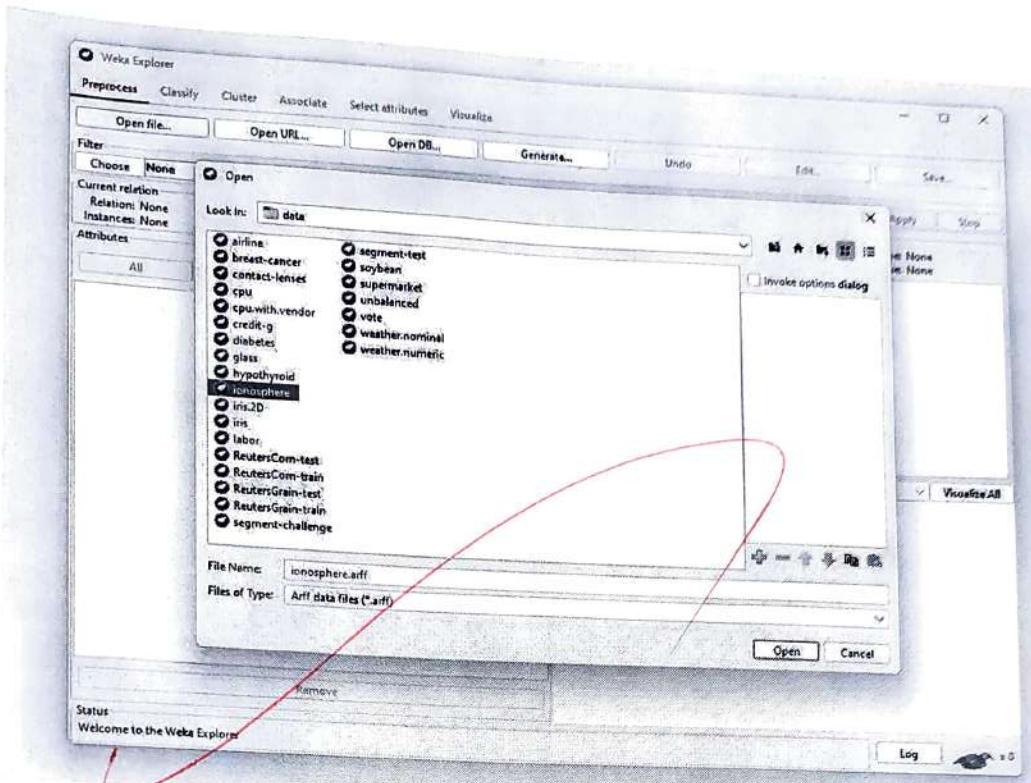
Voting

Ensemble method :- The Voting ensemble technique where multiple ~~of~~ models (classifiers) are combined to improve prediction accuracy. In this approach several models like decision trees, Naive Bayes and SVM are trained on the same dataset and their predictions aggregated by voting for classification. Voting can be either majority voting or weighted voting.

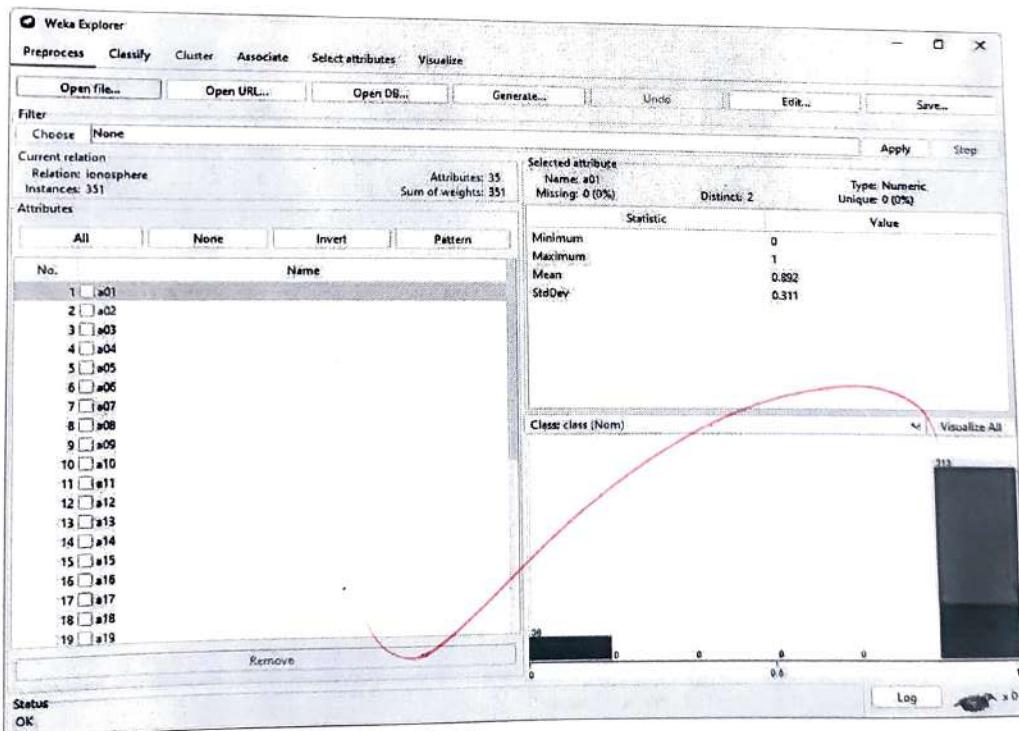
~~Single~~

~~classifier vs Voting ensemble :- Single classifier such as decision trees make predictions based on one model's perspective, on which may lead to bias or variance issue, in contrast the Voting ensemble method combine the strengths of multiple classifier leading to higher accuracy.~~

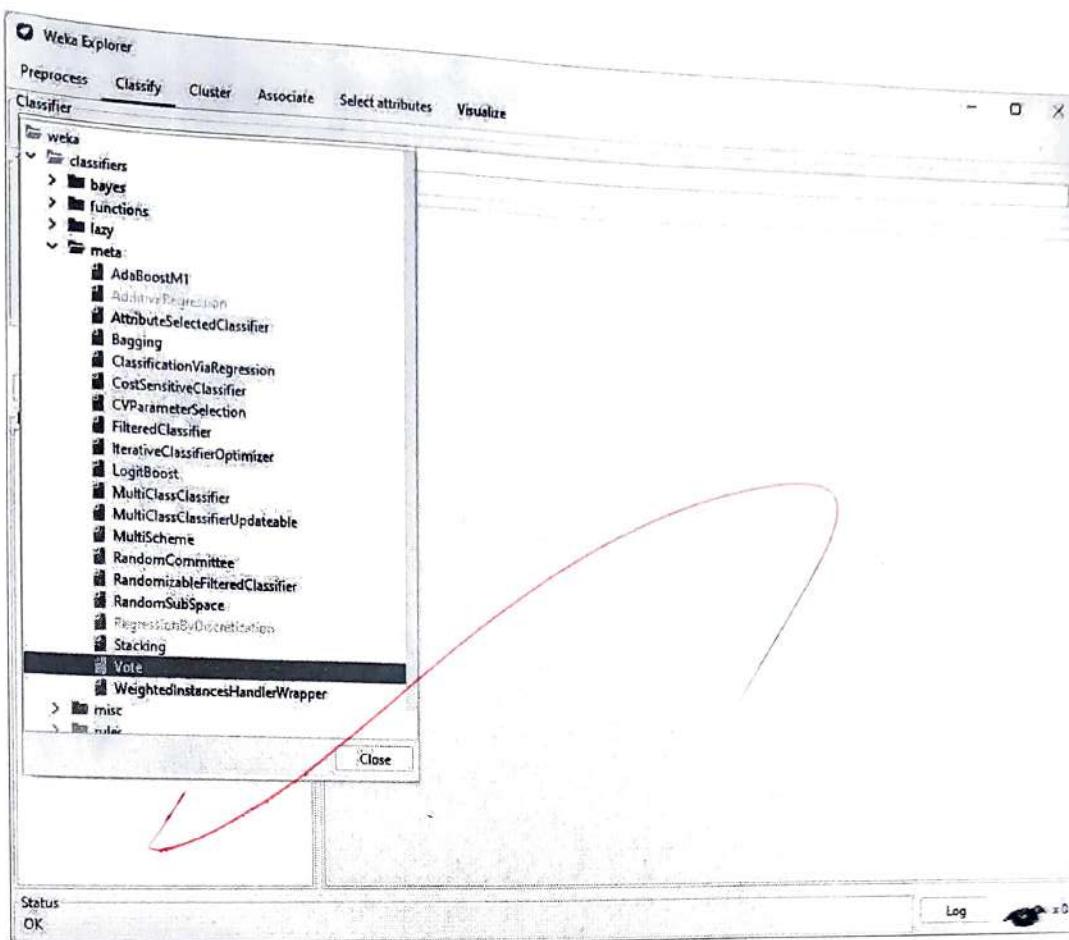
Step 1:



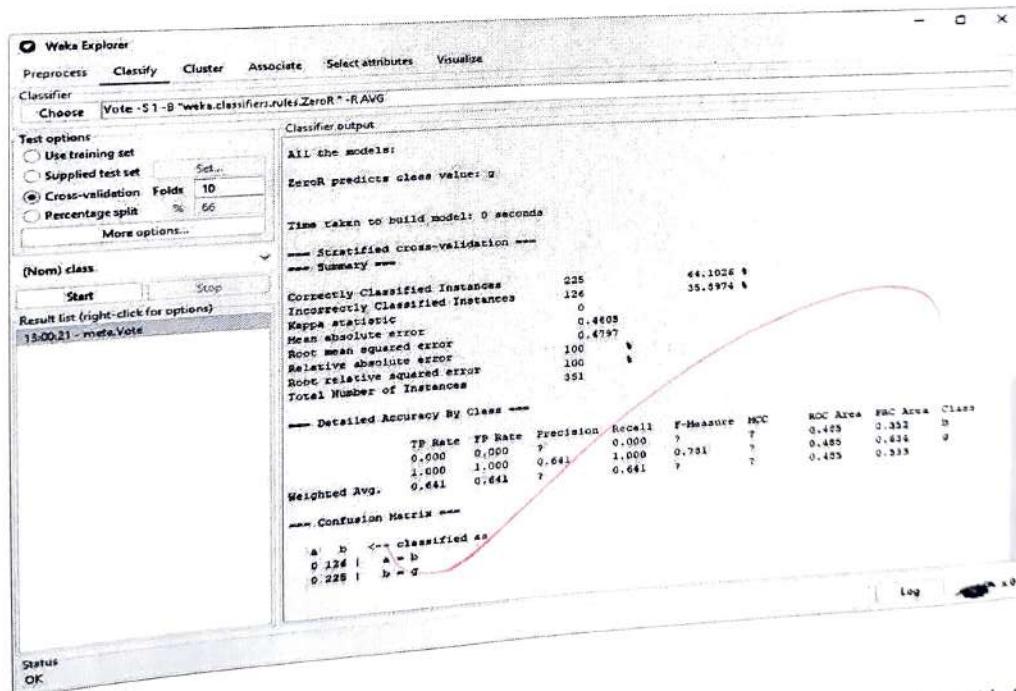
Step 2:



Step 3:

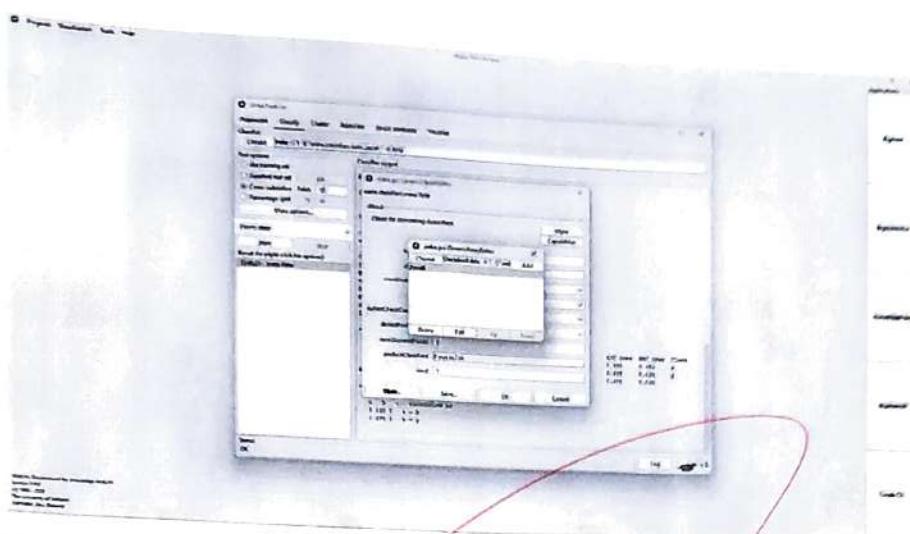


Step 4:

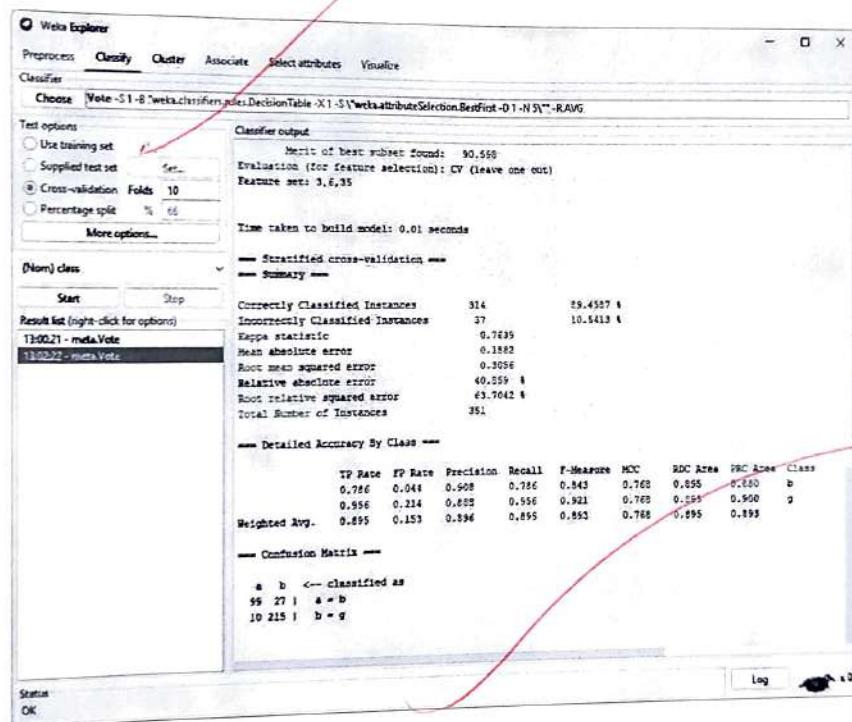


21117711922_HarshitaSharma

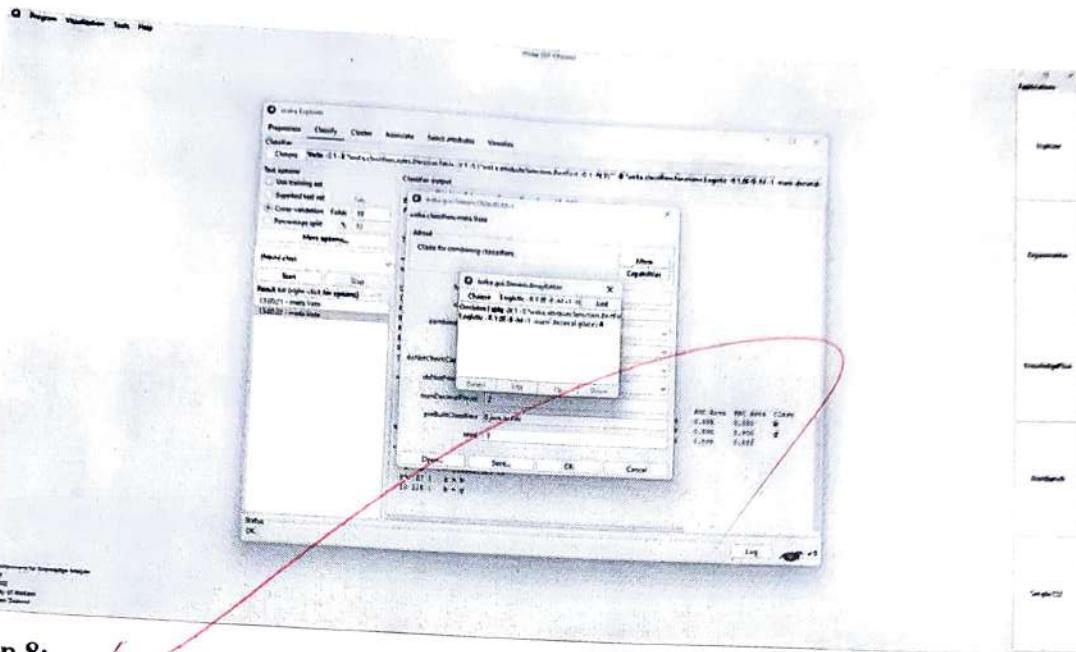
Step 5:



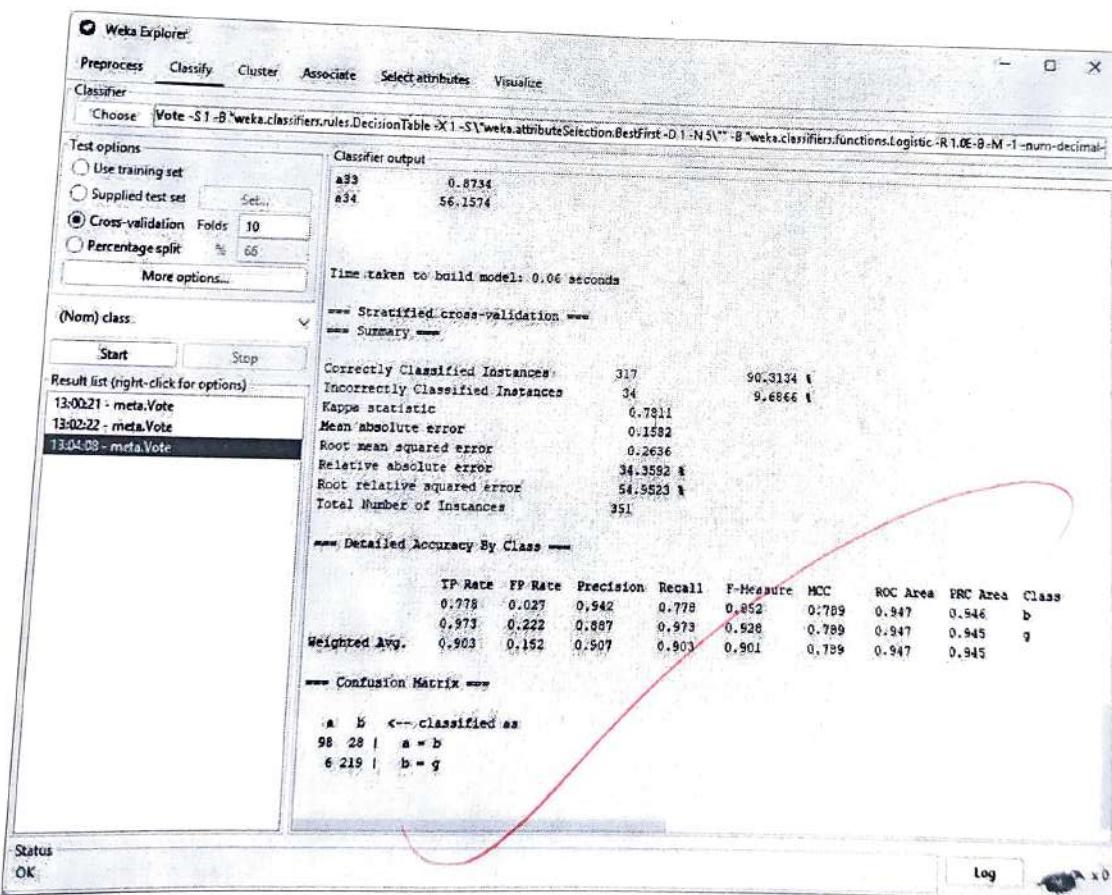
Step 6:



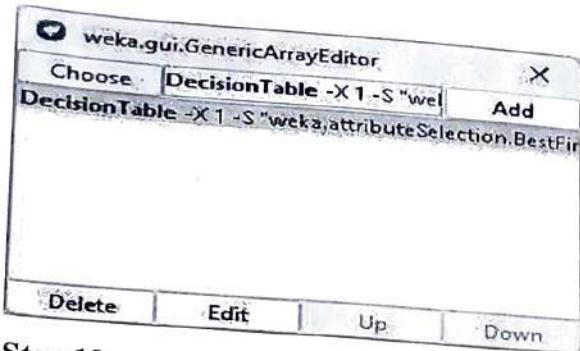
Step 7:



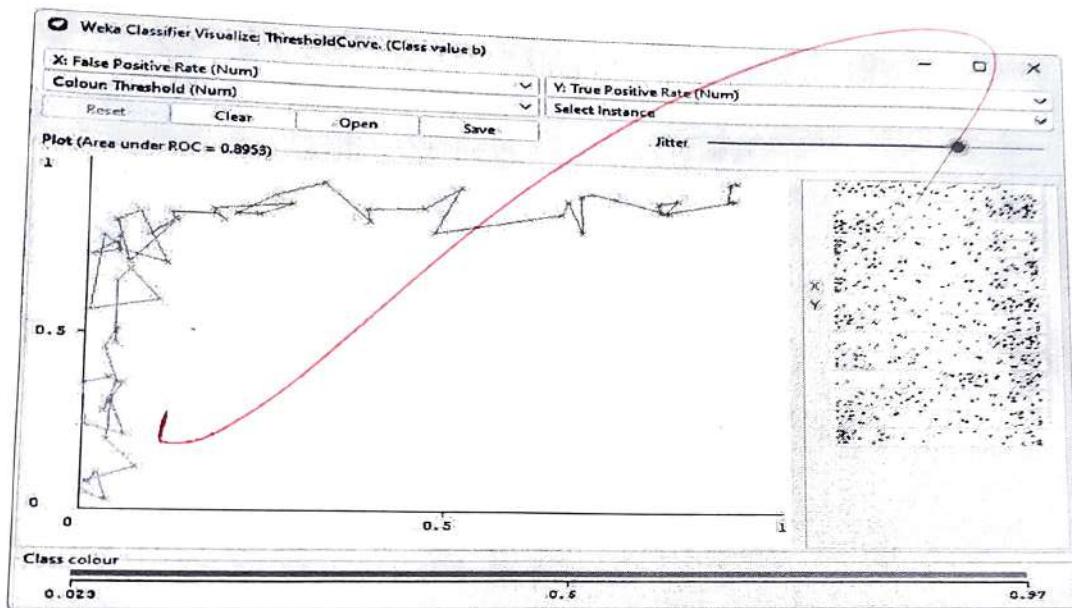
Step 8:



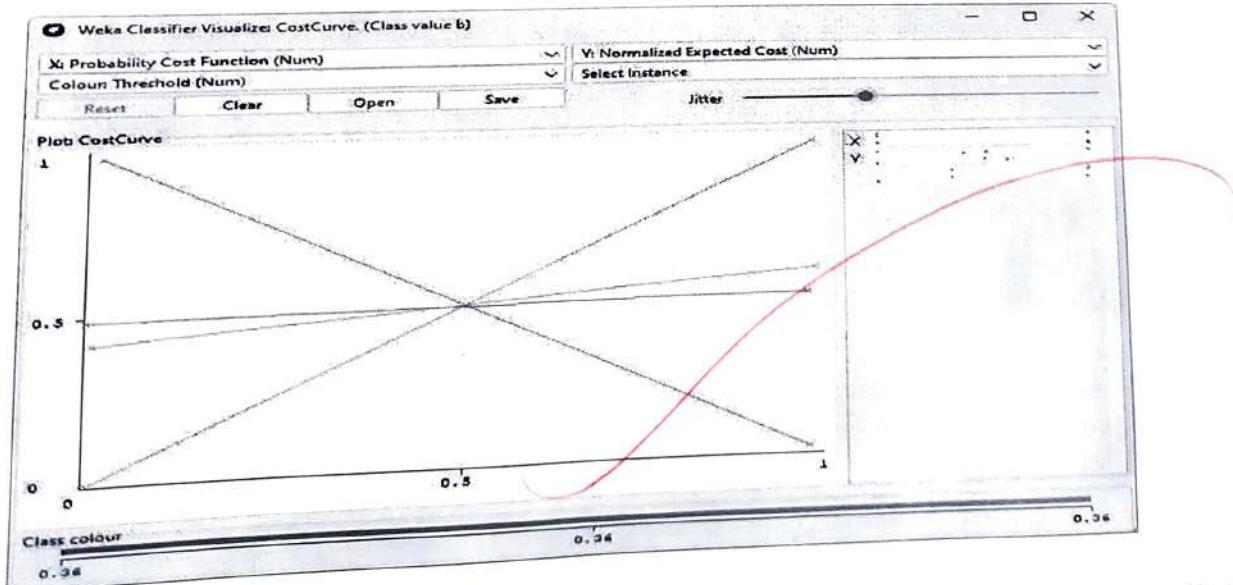
Step 9:



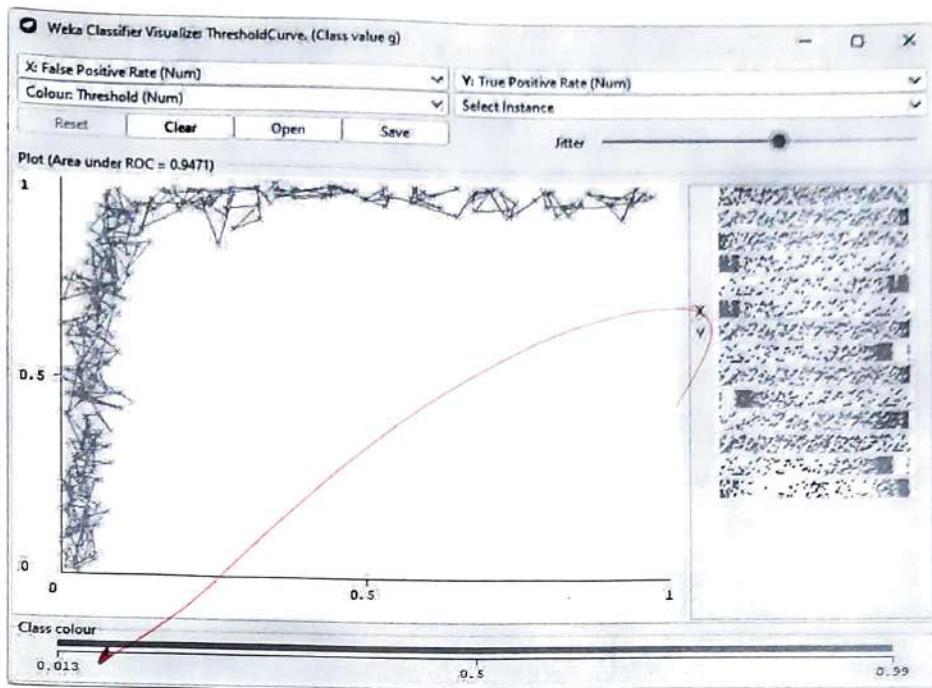
Step 10:



Step 11:



Step 12:



Learning Outcomes:

→ learnt how the voting ensemble method improves accuracy and robustness over single classifiers and gained practical experience in applying and comparing ensemble method with - in same class or voter.

Vote

Experiment 10

Objective: Study of DBMINER tool.

Theory:

DBMiner is an advanced data mining system developed to address the emerging requirements for meaningful knowledge discovery in large relational databases. Since organizations have begun to accumulate enormous amounts of data, traditionally databases designed mainly for storage and query operations are not enough to discover vital insights. DBMiner bridges this gap with a user-friendly yet powerful tool that is capable of multi-dimensional analysis with interactive knowledge discovery. Its primary goal is to enable multi-abstraction-level data exploration allowing organisations to carry out completely transparent analysis of complex datasets. The DBMiner architecture consists of three core components:-

- 1) a relational database that serves as the source of information.
- 2) concept hierarchies so that the user can generalize from one level.
- 3) the knowledge discovery modules, each of which specializes in specific mining tasks.
→ These modules form the backbone functionality of DBMiner, which involves characterization, comparison, classification, association, prediction and clustering.

Advanced data mining techniques boost the capabilities of DBMiner significantly and make it good for intricate analytical tasks. Multiple-level characterization makes use of generalization for showing data at different levels of detail and gives users broad overview before going into details.

The discriminant rules are yet another technique in DBMiner; it finds different groups of data with contrasting features for analysis and determines which features differentiate one group from another, hence revealing the distinguishing features. DBMiner further supports multiple level association rules not only with a single attribute but also among various attributes. for example:- intra-attribute associations in a retail database may point to trends in product categories, whereas, inter-attribute associations may provide you insight in terms of relating purchases with the demographics of customers.

Classification and prediction are two main functionalities in DBMiner tool, more specifically, the classifier module is comprised of combinations between attribute-oriented and decision tree induction methods for classifying the data accordingly, to historical trends. This module is very useful for customer profiling, it can predict which customer segment a new person might belong to according to their attributes.

Another very useful module is prediction, it uses patterns found in the data to predict the future events. For example; how much sales will be during the quarter or which of customers at risk will churn. Clustering groups similar algorithm. DBMiner supports both categorical and numerical data and has the possibility of multi-level clustering, which allows the groupings to be regarded at various levels of detail.

Looking ahead, DBMiner is being further developed in terms of performance to facilitate easier management of increasingly large datasets, and expanded capabilities needs to be made with specialized databases handling spatial, temporal, or multimedia data. In further improving its utility with data-intensive applications, the architecture of DBMiner tool has been designed to support parallel processing with substantially faster computation involving distributed environments and reduced processing time as well as increased stabilities. DBMiner is one of the versatile and robust tool for data mining pertaining to industrial requirements for data-driven insights. It well integrate a relational data-base, concept hierarchies, and knowledge

discovery modules to cater to the interactive exploration of data at different abstraction levels. DBMiner provides actionable patterns and trends by using specialized mining techniques such as → characterization

- association
- classification
- and clustering and the combinations of DMQL (Data Mining Query language) provides more access and usability for users accustomed to SQL.

DBMiner tool epitomizes a robust solution in which organizations can capitalize on their assets relating to data so they can produce strategies, and strategic decision, thereby staying abreast with the rapidly shifting technological landscape in data science and analytics.

Wishful
8/11/24

E-commerce Customer Retention: A Data- Driven Analysis

Presented by
Harshita Sharma
21117711922
ADS-B
G-2

Index

Introduction
Problem Statement
Methodology
Result
Conclusion



v v v v v

Introduction

- In today's competitive landscape, understanding customer behavior is crucial for creating targeted marketing strategies.
- This project leverages data mining techniques, specifically using Weka, to uncover hidden patterns in customer data.
- By analyzing factors like age, payment methods, and product categories, this study aims to segment customers and predict purchasing behaviors, providing actionable insights to drive strategic decisions.
- This project explores customer purchasing behavior using data mining techniques within Weka, a powerful machine learning tool.
- The goal is to analyze customer attributes, segment users by payment preferences, and develop predictive models that help businesses understand purchasing patterns.



Problem Statement

Customer behavior and preferences can be complex and influenced by various factors, such as age, payment methods, and product type. This project seeks to address the challenge of identifying meaningful patterns in customer data to improve business strategies. Specifically, We aim to cluster customers based on their payment methods and create classifiers to predict customer purchasing trends. These insights will help businesses make data-driven decisions for customer engagement, personalized marketing, and inventory management.

Core Challenge: Customer purchasing decisions are influenced by a variety of factors, including demographics, payment methods, and product categories. Understanding these patterns is essential for effective customer segmentation and targeted marketing.

Scope of Study: This project focuses on segmenting customers based on payment methods and creating predictive classifiers for purchasing trends. The data is processed, visualized, and analyzed using Weka, where various clustering and classification techniques provide insights into customer behavior.

Relevance of Study: Analyzing customer purchasing patterns allows businesses to develop strategies for customer retention, personalized marketing, and sales forecasting. This study bridges the gap between raw customer data and actionable business insights.

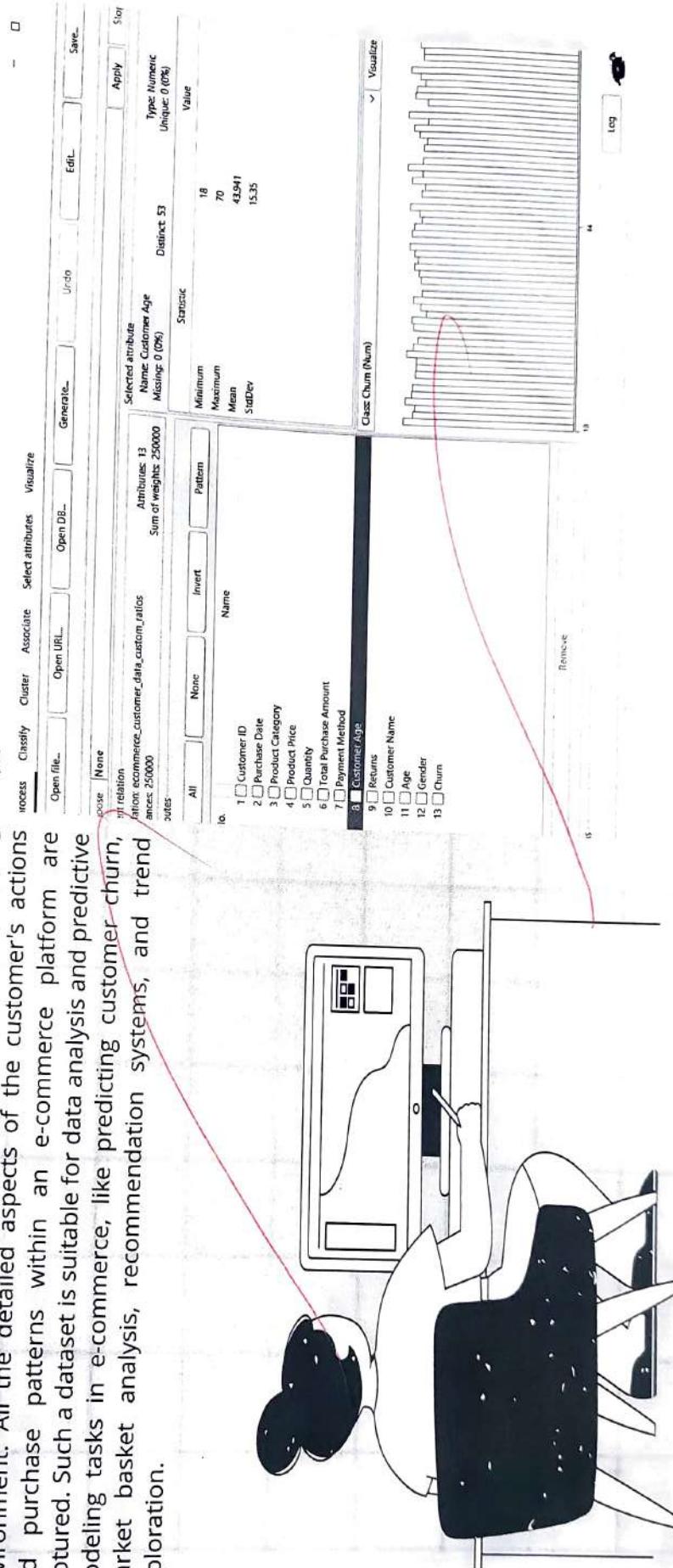
01.

02.

03.

Methodology

This "E-commerce Customer Behavior and Purchase Dataset" is taken from the Kaggle, to simulate a real online shopping environment. All the detailed aspects of the customer's actions and purchase patterns within an e-commerce platform are captured. Such a dataset is suitable for data analysis and predictive modeling tasks in e-commerce, like predicting customer churn, market basket analysis, recommendation systems, and trend exploration.



Pre-Processing Stage

Original Dataset

- Attributes: 13 attributes, including Customer ID, Purchase Date, Product Category, Product Price, Quantity, Total Purchase Amount, Payment Method, Customer Age, Returns, Customer Name, Age, Gender, & Churn.

• **Instances: 250,000 rows.**

• **Selected Attribute:** The Customer Age attribute with values ranging from 18 to 70.

• mean age (43.941)

• standard deviation: (15.35)

After Preprocessing and Noise Removal:

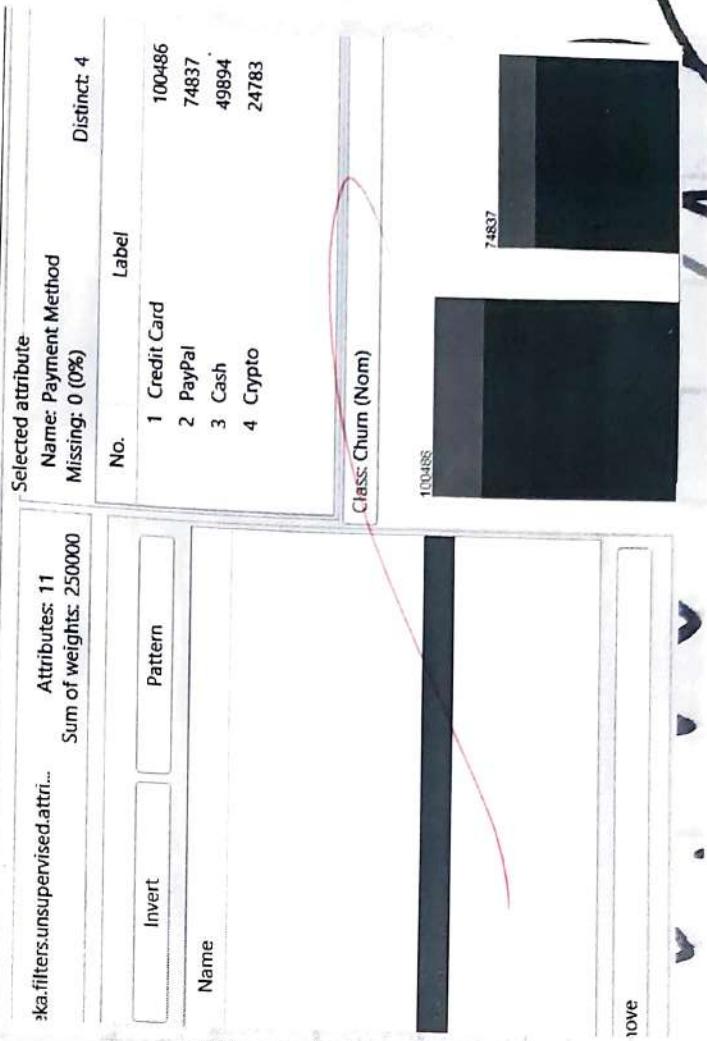
- **Normalization:** The Product Price attribute has been normalized, with a new range from 0 to 1.
 - The mean product price value after normalization is approximately **0.499**, with a standard deviation of **0.289**.
- **Attributes:** processed dataset now has 11 attributes instead of the original 13,
- **Instances:** The number of rows remains the same at **250,000**.

Analysis

- Filter Applied: Discretize.
- Selected Attribute: Payment Method (nominal) with 4 categories:
 - Credit Card: 100,486
 - PayPal: 74,837
 - Cash: 49,894
 - Crypto: 24,783
- Churn Distribution:
 - Bar chart shows churn (red) vs. retention (blue) within each payment method.
 - Credit Card is the most common payment method.
- Insight: Differences in churn across payment methods suggest areas for targeted retention strategies.

Discretization

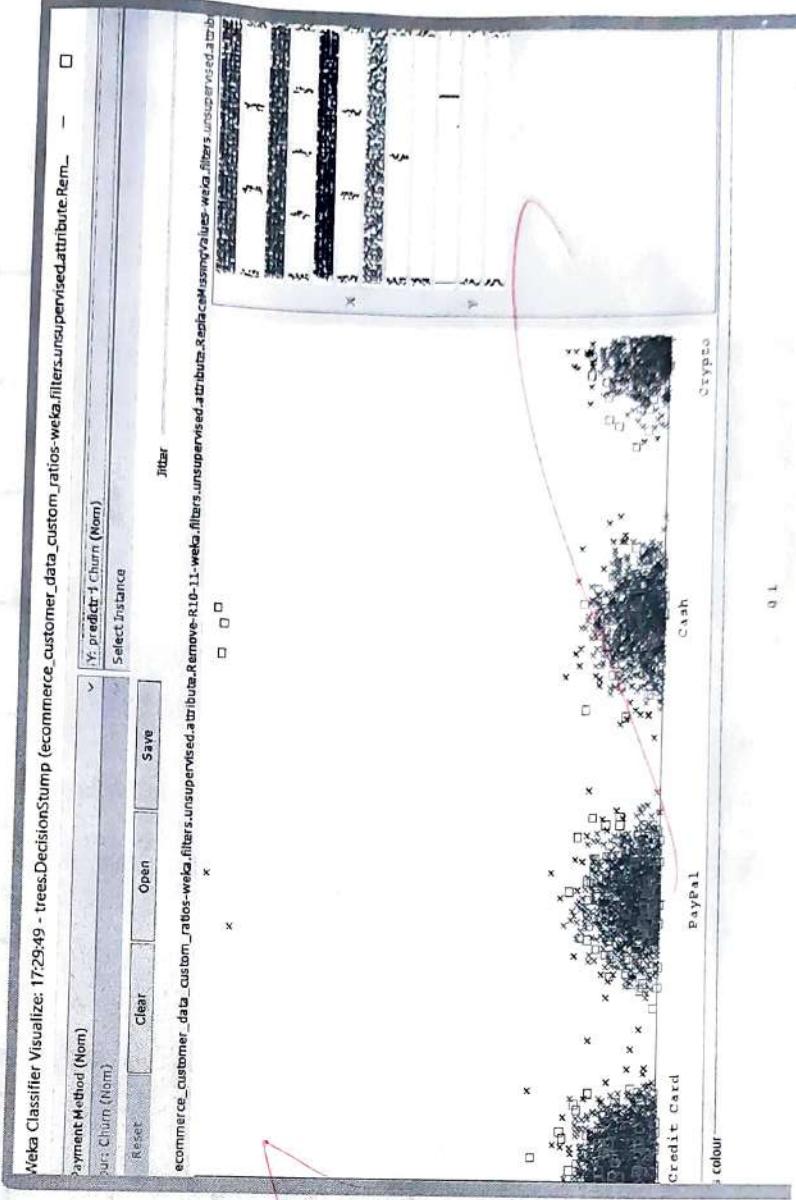
on 6



Decision Tree stump classifier Algorithm

- The plot shows the distribution of customers across different payment methods and their predicted churn status. The color coding helps visualize the accuracy of the model's predictions.

- The algorithm used to generate this decision tree is likely a **Decision Stump**, which is a simple decision tree with only one internal node and two leaf nodes. It is a basic form of a decision tree that can be used for classification tasks like this one.



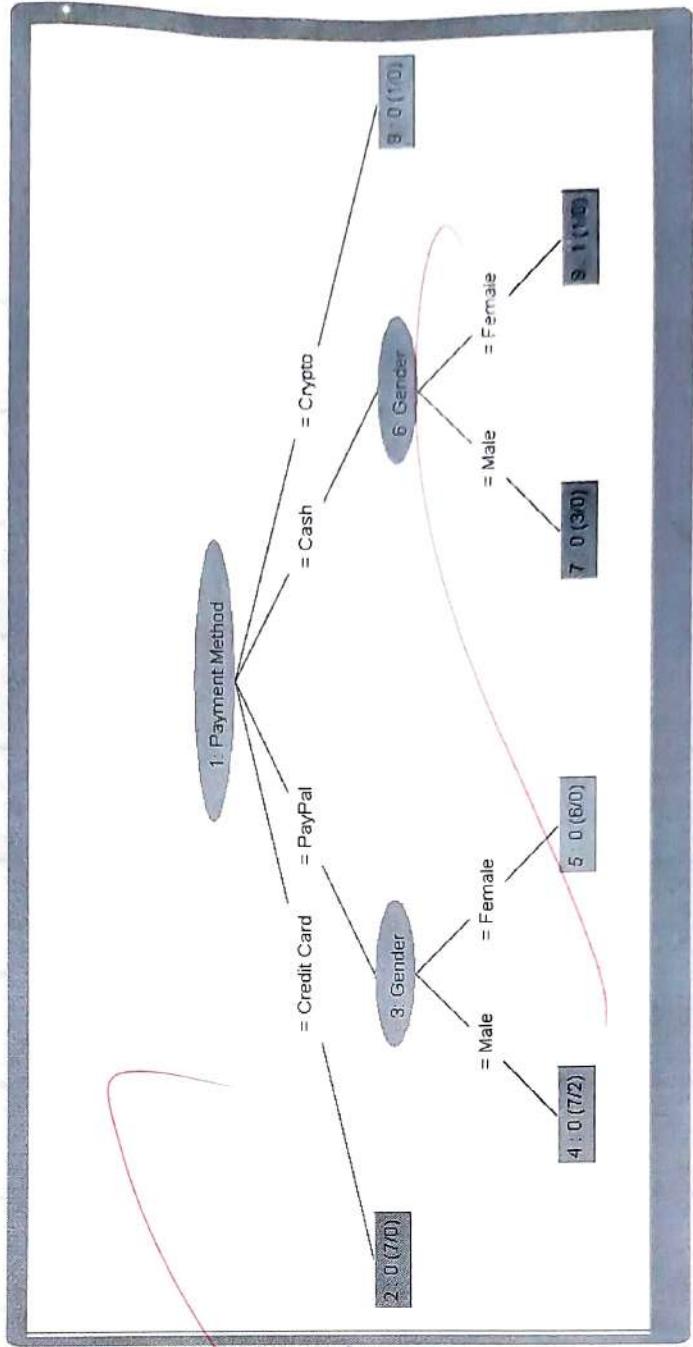
Decision Tree

- Decision trees are a supervised learning algorithm used for both classification and regression tasks.
They work by creating a tree-like model of decisions and their possible consequences.

- **Root Node:** Splits data based on "Payment Method".

- **Child Nodes:** Further split based on "Gender".

- **Leaf Nodes:** Represent final prediction (churn or not churn) with associated instance counts



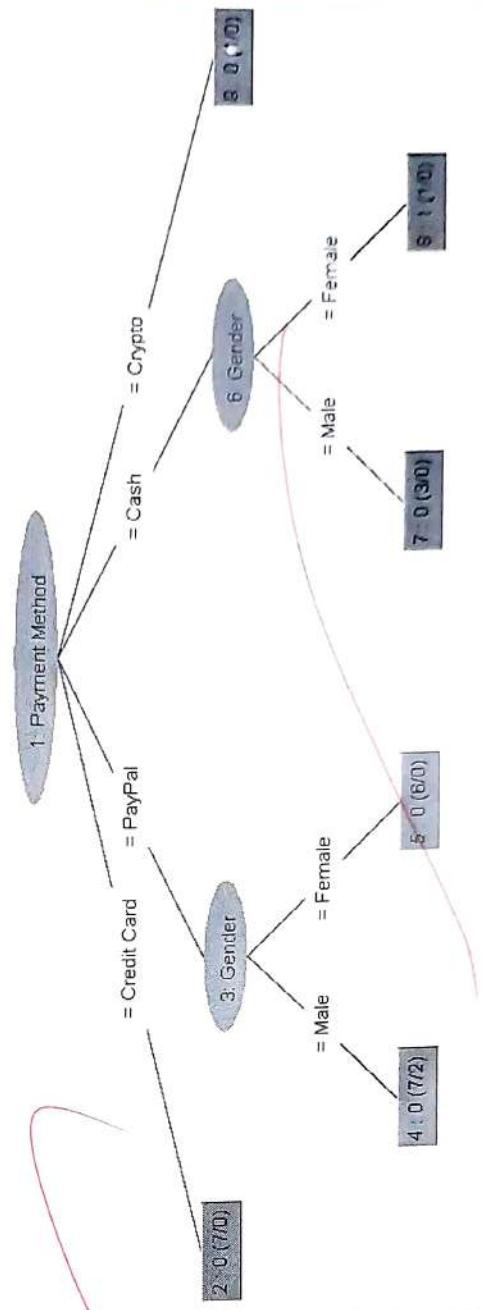
Decision Tree

- Decision trees are a supervised learning algorithm used for both classification and regression tasks.
They work by creating a tree-like model of decisions and their possible consequences.

• **Root Node:** Splits data based on "Payment Method".

• **Child Nodes:** Further split based on "Gender".

• **Leaf Nodes:** Represent final prediction (churn or not churn) with associated instance counts

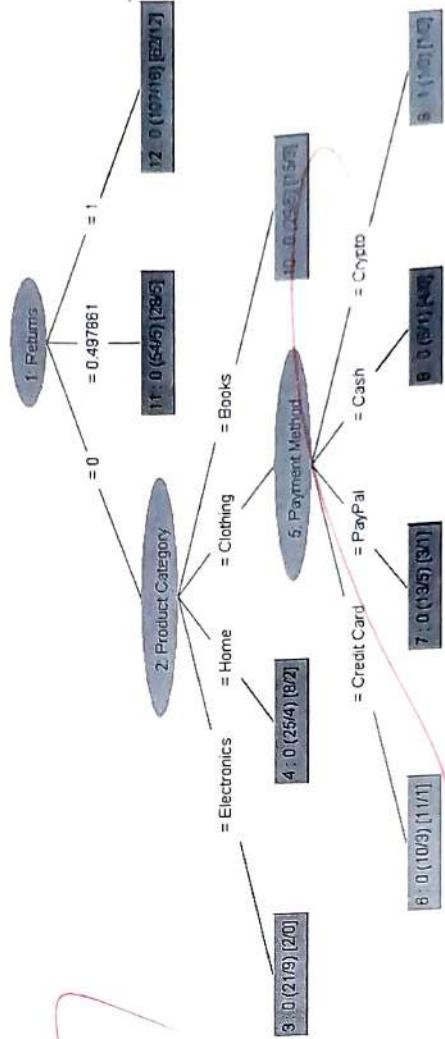


REPTREE

Analysis

Weka Classifier Tree Visualizer: 16:57:00 - trees.REPTree (ecommerce_customer_data_custom_ratios-weka.filters.unsupervised.attribute.Remove-R10-11
Tree View

- REPTree (Reduced Error Pruning Tree) is a decision tree algorithm used in data mining for classification and regression tasks. It's known for its efficiency and simplicity.
- This tree is likely used to predict customer churn based on various factors like product returns, product category, and payment method.



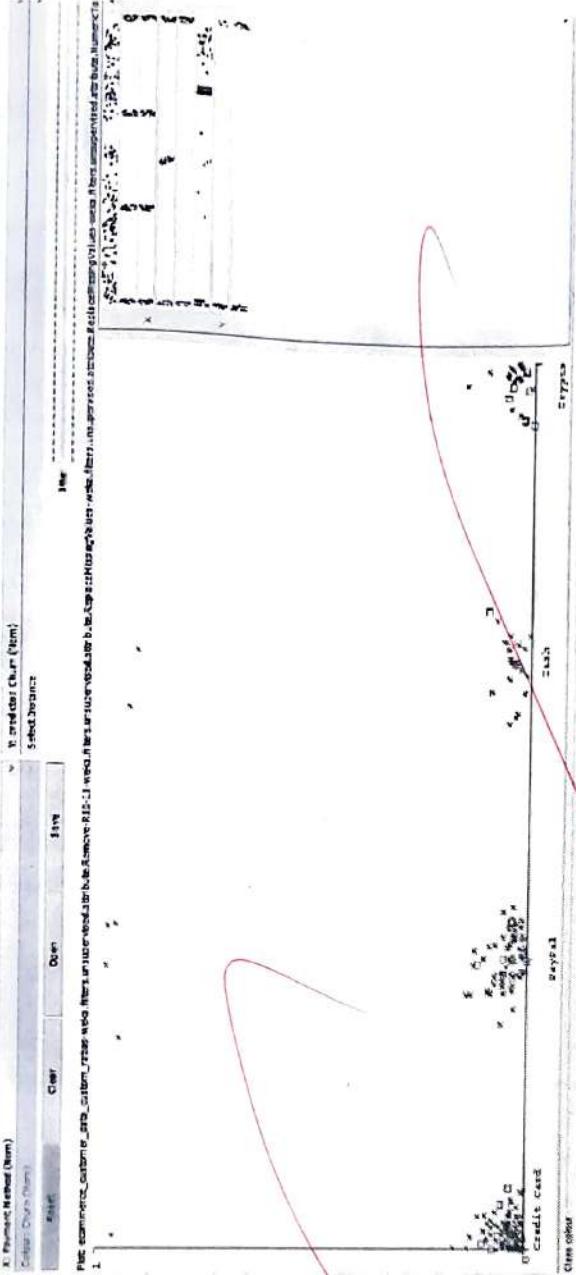
Bagging

Explanation:

- The plot displays the performance improvement achieved by bagging, which trains multiple models on bootstrapped samples of the data and combines them to form a more stable, accurate classifier.
- Bagging reduces variance by averaging predictions from several models, which makes it particularly effective when individual models are prone to high variance.

Algorithm Insight:

- By averaging results from diverse models, bagging creates a more reliable and accurate predictor. It is especially useful in datasets where individual classifiers may overfit.
- Bagging works well with algorithms like decision trees, which are inherently high-variance models. It increases stability without significantly increasing bias.



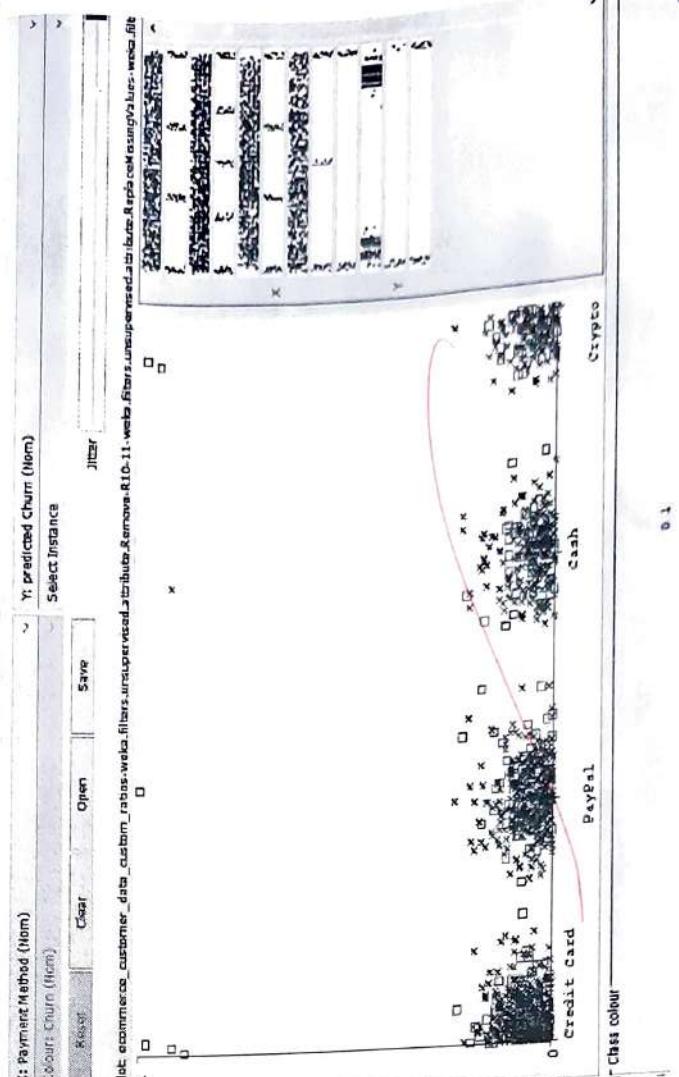
AdaBoostM1 Algorithm

Explanation:

- The plot illustrates the error rate of the AdaBoostM1 algorithm over several boosting rounds, showing how misclassified samples are given more focus in each iteration.
- AdaBoostM1 is an ensemble method that improves accuracy by iteratively adjusting the weights of instances, especially those that are misclassified, to "boost" the predictive power of the overall model.

Algorithm Insight:

- This model is effective for handling difficult-to-classify instances by combining weak learners into a strong predictive model. It is particularly suited for cases where improving accuracy on challenging samples is crucial.
- AdaBoostM1 is sensitive to noisy data, as misclassified noisy instances may receive higher weights, which can sometimes lead to overfitting.



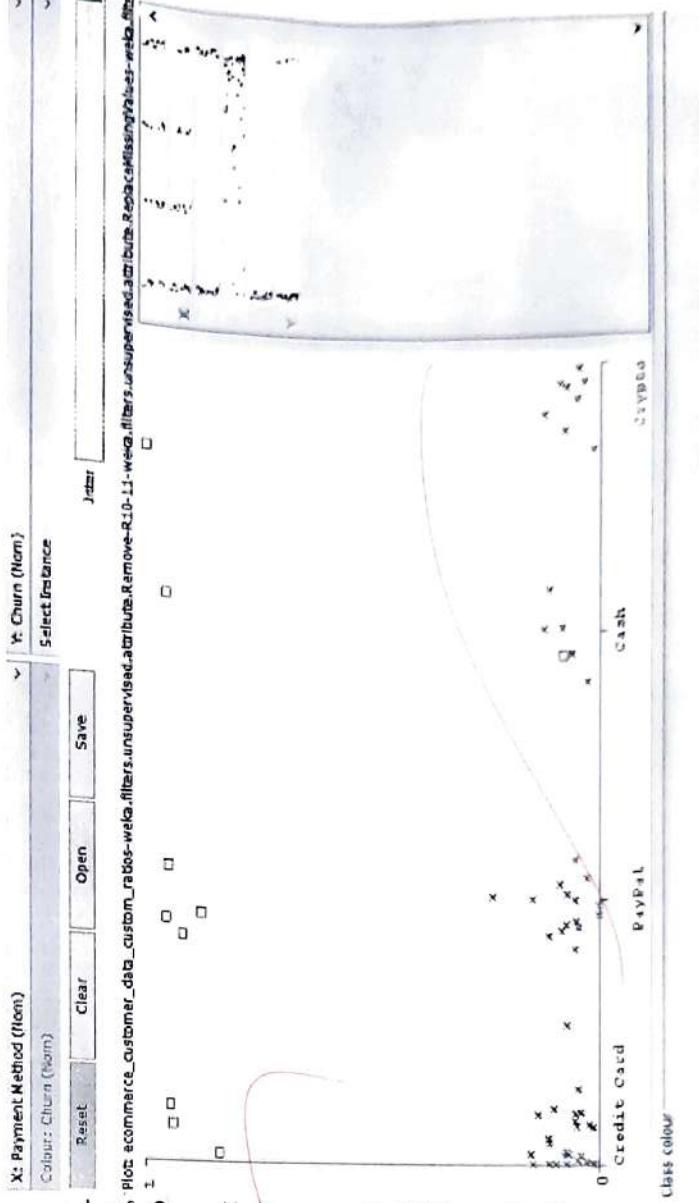
Meta-Vote Ensemble Algorithm

Explanation:

- The plot shows the classification accuracy of the Meta-Vote ensemble, which combines J48, Random Forest, and Naïve Bayes classifiers by majority voting to enhance predictive performance.
- Meta-Vote ensembles leverage the strengths of different algorithms to achieve a higher overall accuracy by aggregating their individual predictions.

Algorithm Insight:

- This approach benefits from the complementary strengths of various models. While J48 provides interpretability, Random Forest offers robustness, and Naïve Bayes adds probabilistic reasoning.
- Meta-Vote ensembles are particularly useful when individual models have diverse strengths that, when combined, provide a balanced prediction.



Simple Logistic Algorithm

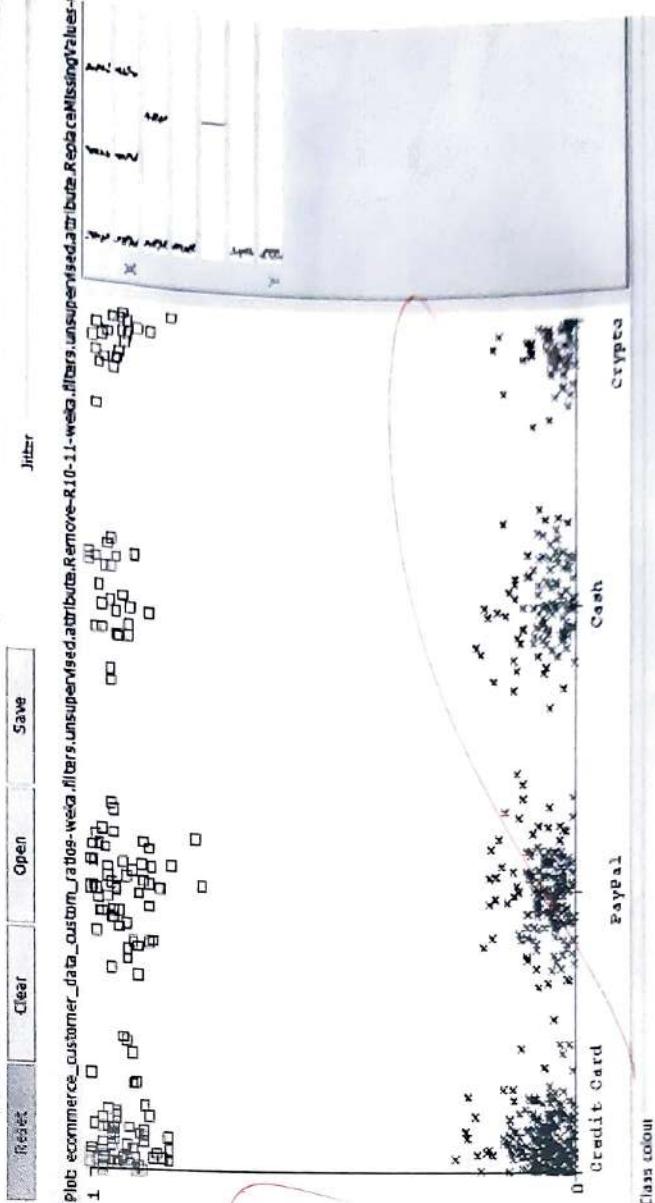
Explanation:

The plot illustrates the decision boundaries created by the Simple Logistic algorithm, which predicts probabilities for each class and is ideal for linear classification.

Simple Logistic is a regression-based classifier that provides probabilities, offering a straightforward interpretation for binary or multi-class classification.

Algorithm Insight:

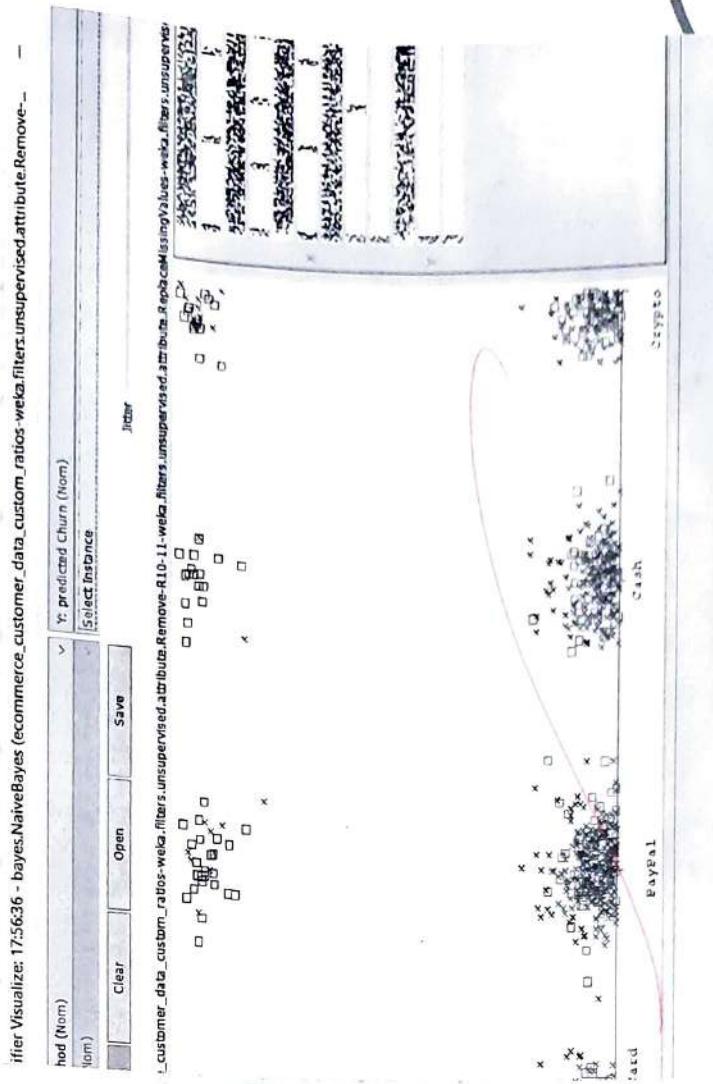
- It is best suited for problems with linearly separable classes and offers a probabilistic output, which is useful for decision-making.
- While limited to linear relationships, it is a simple and interpretable model, effective for datasets with binary or few categories.



Naïve Bayes

Analysis

- Naïve Bayes is a probabilistic classification algorithm based on Bayes' theorem.
- It assumes that the features are conditionally independent,
- Meaning that the occurrence of one feature does not influence the probability of another feature given the class.

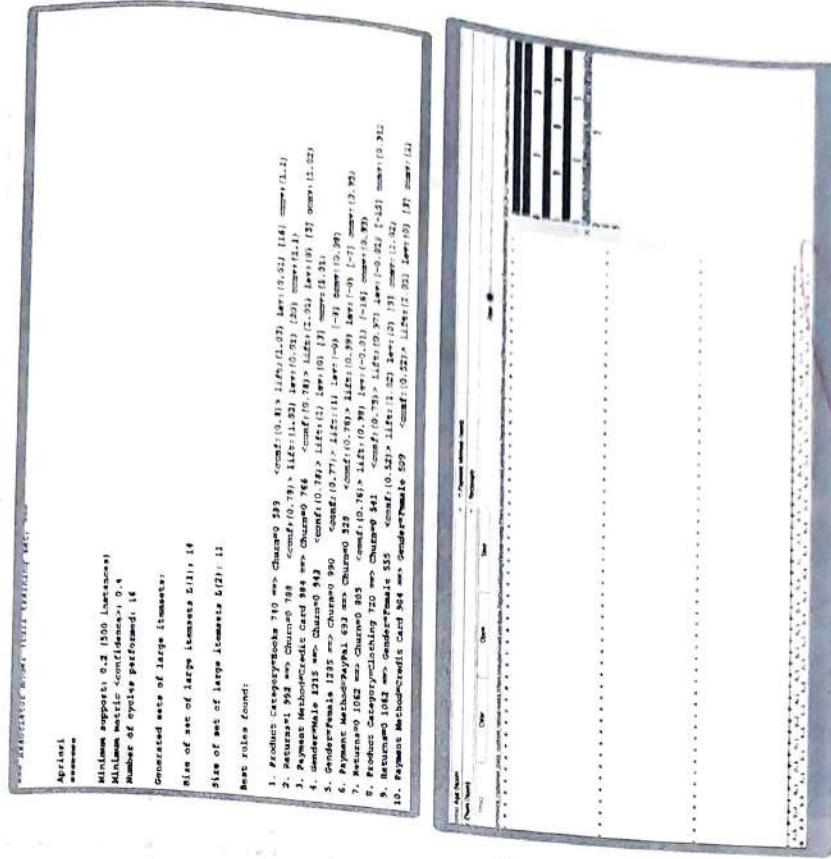


Apriori Algorithm

- The first image shows the results of an Apriori association rule mining algorithm applied to the e-commerce dataset. This analysis reveals interesting patterns in customer behavior and product relationships.
- The second image presents a scatter plot visualization of customer data, focusing on customer age and payment method.

Insights

- Customer Churn:** Certain customer demographics, product returns, and payment methods are linked to higher churn rates.
- Product Relationships:** Specific products are frequently purchased together, offering opportunities for cross-selling and upselling.
- Age and Payment Method:** Different age groups prefer different payment methods.
- Targeted Marketing:** Tailor marketing campaigns to specific age groups and customer segments.
- Churn Prevention:** Implement targeted retention strategies for at-risk customers.
- Product Recommendations:** Leverage association rules for personalized product recommendations.
- Payment Optimization:** Offer a variety of secure payment options to cater to different preferences.



Result

Our analysis yielded meaningful insights into customer segments:

- Clustering Outcome: Customer groups based on payment methods provided clear visual insights. For example, Orange represents active use of payment methods by age group, while Blue indicates inactivity.
- Classification Performance: The ensemble models, particularly AdaBoostM1 and meta-vote ensemble, showed the best accuracy, with error plots highlighting performance across classifiers.
- Error Analysis: Decision stump, Naïve Bayes, and other models displayed varying error rates, helping identify areas for potential improvement.



Conclusion

- In conclusion, this project highlights the significant value of data mining for customer analysis, using Weka to facilitate in-depth pattern recognition and predictive modeling.
- Through the integration of clustering and classification techniques, we identified key behavioral trends that can be strategically leveraged to improve customer engagement and drive business growth.
- The K-means clustering revealed distinct customer segments based on payment methods, providing actionable insights into age-based preferences and purchasing behaviors.
- Additionally, the use of ensemble methods, particularly the meta-vote and AdaBoostM1, offered high accuracy in classification, underscoring the robustness of combining models for enhanced prediction.
- These insights equip businesses to make data-driven decisions, personalize marketing efforts, and optimize resources to better align with customer needs.



Future Work

- This project underscores how data mining is not only a tool for analysis but also a pathway to strategic advancements in customer satisfaction and retention.
- Future work can expand upon this foundation by incorporating larger datasets and refining models to achieve even greater predictive precision, thereby continually enhancing the customer experience.



**THANK
YOU**

