

Degree Engineering

A Laboratory Manual for

Data Mining

(3160714)

B.E. Semester 6
(Computer Engineering)



**Directorate of Technical Education, Gandhinagar,
Gujarat**

Preface

Main motto of any laboratory/practical/field work is for enhancing required skills as well as creating ability amongst students to solve real time problem by developing relevant competencies in psychomotor domain. By keeping in view, GTU has designed competency focused outcome-based curriculum for engineering degree programs where sufficient weightage is given to practical work. It shows importance of enhancement of skills amongst the students and it pays attention to utilize every second of time allotted for practical amongst students, instructors and faculty members to achieve relevant outcomes by performing the experiments rather than having merely study type experiments. It is must for effective implementation of competency focused outcome-based curriculum that every practical is keenly designed to serve as a tool to develop and enhance relevant competency required by the various industry among every student. These psychomotor skills are very difficult to develop through traditional chalk and board content delivery method in the classroom. Accordingly, this lab manual is designed to focus on the industry defined relevant outcomes, rather than old practice of conducting practical to prove concept and theory.

By using this lab manual students can go through the relevant theory and procedure in advance before the actual performance which creates an interest and students can have basic idea prior to performance. This in turn enhances pre-determined outcomes amongst students. Each experiment in this manual begins with competency, industry relevant skills, course outcomes as well as practical outcomes (objectives). The students will also achieve safety and necessary precautions to be taken while performing practical.

This manual also provides guidelines to faculty members to facilitate student centric lab activities through each experiment by arranging and managing necessary resources in order that the students follow the procedures with required safety and necessary precautions to achieve the outcomes. It also gives an idea that how students will be assessed by providing rubrics.

Data mining is a key to sentiment analysis, price optimization, database marketing, credit risk management, training and support, fraud detection, healthcare and medical diagnoses, risk assessment, recommendation systems and much more. It can be an effective tool in just about any industry, including retail, wholesale distribution, service industries, telecom, communications, insurance, education, manufacturing, healthcare, banking, science, engineering, and online marketing or social media.

Utmost care has been taken while preparing this lab manual however always there is chances of improvement. Therefore, we welcome constructive suggestions for improvement and removal of errors if any.

Vishwakarma Government Engineering College
Department of Computer Engineering

CERTIFICATE

This is to certify that Mr./Ms. Abhi Yogeshumar Patel Enrollment No 220170107076 of B.E. Semester 6th from **Computer Engineering Department** of this Institute (GTU Code: 017) has satisfactorily completed the Practical / Tutorial work for the subject **Data Mining (3160714)** for the academic year 2024-25.

Place: _____

Date: _____

Signature of Course Faculty

Head of the Department

DTE's Vision

- To provide globally competitive technical education
- Remove geographical imbalances and inconsistencies
- Develop student friendly resources with a special focus on girls' education and support to weaker sections
- Develop programs relevant to industry and create a vibrant pool of technical professionals

Institute's Vision

- To create an ecosystem for proliferation of socially responsible and technically sound engineers, innovators and entrepreneurs.

Institute's Mission

- To develop state-of-the-art laboratories and well-equipped academic infrastructure.
- To motivate faculty and staff for qualification up-gradation, and enhancement of subject knowledge.
- To promote research, innovation and real-life problem-solving skills.
- To strengthen linkages with industries, academic and research organizations.
- To reinforce concern for sustainability, natural resource conservation and social responsibility.

Department's Vision

- To create an environment for providing value-based education in Computer Engineering through innovation, team work and ethical practices.

Department's Mission

- To produce computer engineering graduates according to the needs of industry, government, society and scientific community.
- To develop state of the art computing facilities and academic infrastructure.
- To develop partnership with industries, government agencies and R & D organizations for knowledge sharing and overall development of faculties and students.
- To solve industrial, governance and societal issues by applying computing techniques.
- To create environment for research and entrepreneurship.

Programme Outcomes (POs)

1. **Engineering knowledge:** Apply the knowledge of mathematics, science, engineering fundamentals, and an engineering specialization to the solution of complex engineering problems.
2. **Problem analysis:** Identify, formulate, review research literature, and analyze complex engineering problems reaching substantiated conclusions using first principles of mathematics, natural sciences, and engineering sciences.
3. **Design/development of solutions:** Design solutions for complex engineering problems and design system components or processes that meet the specified needs with appropriate consideration for the public health and safety, and the cultural, societal, and environmental considerations.
4. **Conduct investigations of complex problems:** Use research-based knowledge and research methods including design of experiments, analysis and interpretation of data, and synthesis of the information to provide valid conclusions.
5. **Modern tool usage:** Create, select, and apply appropriate techniques, resources, and modern engineering and IT tools including prediction and modeling to complex engineering activities with an understanding of the limitations.
6. **The engineer and society:** Apply reasoning informed by the contextual knowledge to assess societal, health, safety, legal and cultural issues and the consequent responsibilities relevant to the professional engineering practice.
7. **Environment and sustainability:** Understand the impact of the professional engineering solutions in societal and environmental contexts, and demonstrate the knowledge of, and need for sustainable development.
8. **Ethics:** Apply ethical principles and commit to professional ethics and responsibilities and norms of the engineering practice.
9. **Individual and team work:** Function effectively as an individual, and as a member or leader in diverse teams, and in multidisciplinary settings.
10. **Communication:** Communicate effectively on complex engineering activities with the engineering community and with society at large, such as, being able to comprehend and write effective reports and design documentation, make effective presentations, and give and receive clear instructions.
11. **Project management and finance:** Demonstrate knowledge and understanding of the engineering and management principles and apply these to one's own work, as a member and leader in a team, to manage projects and in multidisciplinary environments.
12. **Life-long learning:** Recognize the need for, and have the preparation and ability to engage in independent and life-long learning in the broadest context of technological change.

Program Specific Outcomes (PSOs)

- Sound knowledge of fundamentals of computer science and engineering including software and hardware.
- Develop the software using sound software engineering principles having web based/mobile based interface.
- Use various tools and technology supporting modern software frameworks for solving problems having large volume of data in the domain of data science and machine learning.

Program Educational Objectives (PEOs)

- Possess technical competence in solving real life problems related to Computing.
- Acquire good analysis, design, development, implementation and testing skills to formulate simple computing solutions to the business and societal needs.
- Provide requisite skills to pursue entrepreneurship, higher studies, research, and development and imbibe high degree of professionalism in the fields of computing.
- Embrace life-long learning and remain continuously employable.
- Work and excel in a highly competence supportive, multicultural and professional environment which abiding to the legal and ethical responsibilities.

Practical – Course Outcome matrix**Course Outcomes (COs):**

CO_3160714.1	Perform the preprocessing of data and apply mining techniques on it.
CO_3160714.2	Identify the association rules, classification, and clusters in large data sets.
CO_3160714.3	Solve real-world problems in business and scientific information using data mining.
CO_3160714.4	Use data analysis tools for scientific applications.
CO_3160714.5	Implement various supervised machine learning algorithms.

Sr. No.	Objective(s) of Experiment	CO 1	CO 2	CO 3	CO 4	CO 5
1.	Identify how data mining is an interdisciplinary field by an Application.			√		
2.	<p>Write programs to perform the following tasks of preprocessing (any language).</p> <ul style="list-style-type: none"> 2.1 Noisy data handling <ul style="list-style-type: none"> ➤ Equal Width Binning ➤ Equal Frequency/Depth Binning 2.2 Normalization Techniques <ul style="list-style-type: none"> ➤ Min max normalization ➤ Z score normalization ➤ Decimal scaling 2.3. Implement data dispersion measure Five Number Summary generate box plot using python libraries 					
3.	To perform hand on experiments of data preprocessing with sample data on Orange tool.	√			√	
4.	Implement Apriori algorithm of association rule data mining technique in any Programming language.		√			
5.	Apply association rule data mining technique on sample data sets using WEKA.		√		√	
6.	Apply Classification data mining technique on sample data sets in WEKA.				√	√
7.	<p>7.1. Implement Classification technique with quality Measures in any Programming language.</p> <p>7.2 Implement Regression technique in any Programming language.</p>					√
8.	Apply K-means Clustering Algorithm any Programming language.		√		√	
9.	Perform hands on experiment on any advance mining Techniques Using Appropriate Tool.				√	
10.	Solve Real world problem using Data Mining Techniques using Python Programming Language.			√		

Guidelines for Faculty members

1. Teacher should provide the guideline with demonstration of practical to the students with all features.
2. Teacher shall explain basic concepts/theory related to the experiment to the students before starting of each practical
3. Involve all the students in performance of each experiment.
4. Teacher is expected to share the skills and competencies to be developed in the students and ensure that the respective skills and competencies are developed in the students after the completion of the experimentation.
5. Teachers should give opportunity to students for hands-on experience after the demonstration.
6. Teacher may provide additional knowledge and skills to the students even though not covered in the manual but are expected from the students by concerned industry.
7. Give practical assignment and assess the performance of students based on task assigned to check whether it is as per the instructions or not.
8. Teacher is expected to refer complete curriculum of the course and follow the guidelines for implementation.

Instructions for Students

1. Students are expected to carefully listen to all the theory classes delivered by the faculty members and understand the COs, content of the course, teaching and examination scheme, skill set to be developed etc.
2. Students will have to perform experiments as per practical list given.
3. Students have to show output of each program in their practical file.
4. Students are instructed to submit practical list as per given sample list shown on next page.
5. Student should develop a habit of submitting the experimentation work as per the schedule and s/he should be well prepared for the same.

Common Safety Instructions

Students are expected to

- 1) switch on the PC carefully (not to use wet hands)
- 2) shutdown the PC properly at the end of your Lab
- 3) carefully handle the peripherals (Mouse, Keyboard, Network cable etc)
- 4) use Laptop in lab after getting permission from Teacher

Index

(Progressive Assessment Sheet)

Sr. No.	Objective(s) of Experiment	Page No.	Date of performance	Date of submission	Assessment Marks	Sign. of Teacher with date	Remarks
1	Identify how data mining is an interdisciplinary field by an Application.						
2	Write programs to perform the following tasks of preprocessing (any language). <ul style="list-style-type: none"> 2.1 Noisy data handling <ul style="list-style-type: none"> ➤ Equal Width Binning ➤ Equal Frequency/Depth Binning 2.2 Normalization Techniques <ul style="list-style-type: none"> ➤ Min max normalization ➤ Z score normalization ➤ Decimal scaling 2.3. Implement data dispersion measure Five Number Summary generate box plot using python libraries. 						
3	To perform hand on experiments of data preprocessing with sample data on Orange tool.						
4	Implement Apriori algorithm of association rule data mining technique in any Programming language.						
5	Apply association rule data mining technique on sample data sets using WEKA.						
6	Apply Classification data mining technique on sample data sets in WEKA .						
7	7.1. Implement Classification technique with quality Measures in any Programming language. 7.2. Implement Regression technique in any Programming language.						
8	Apply K-means Clustering Algorithm any Programming language.						
9	Perform hands on experiment on any advance mining Techniques Using Appropriate Tool.						
10	Solve Real world problem using Data Mining Techniques using Python Programming Language.						

Experiment No - 1

Aim: Identify how data mining is an interdisciplinary field by an Application.

Data mining is an interdisciplinary field that involves computer science, statistics, mathematics, and domain-specific knowledge. One application that showcases the interdisciplinary nature of data mining

Date: // Write date of experiment here

Competency and Practical Skills: Understanding and Analyzing

Relevant CO: CO3

Objectives:

- (a) To understand the application of domain
- (b) To understand Preprocessing Techniques.
- (c) To understand the application's use of Data Mining functionalities.

Equipment/Instruments: Personal Computer

Theory:

System Name: Image Captioning (Chest X-Rays)

Image captioning for Chest X-Rays is a great example of how data mining is an interdisciplinary field involving multiple expertise areas. The process of creating an image captioning system for medical images involves the following steps:

Dataset:

An image captioning system for Chest X-Rays requires a dataset containing X-ray images and their corresponding textual descriptions. Here are some examples of datasets:

- **Kaggle Chest X-Ray Dataset:** This dataset contains labeled chest X-ray images with corresponding reports, widely used for training and testing deep learning models.
- **MIMIC-CXR:** A large-scale dataset consisting of over 377,000 chest X-rays and corresponding radiology reports, released by MIT for research purposes.
- **NIH Chest X-ray Dataset:** A dataset of over 100,000 chest X-ray images from NIH Clinical Center, labeled with different thoracic diseases.

Preprocessing:

It involves cleaning and transforming the data to make it suitable for analysis. Here are some preprocessing techniques commonly used in medical image captioning systems:

- **Data Cleaning:** Removing poor-quality images, ensuring proper annotations, and handling missing data in radiology reports.
- **Image Normalization:** Adjusting pixel intensity values to a common scale, such as converting all X-ray images to grayscale and normalizing pixel values between 0 and 1.
- **Text Preprocessing:** Cleaning textual reports by removing unnecessary characters, tokenizing sentences, and standardizing medical terminology.
- **Feature Extraction:** Extracting key features from images using convolutional neural networks (CNNs) or other deep learning techniques.
- **Data Augmentation:** Applying transformations such as rotation, zooming, and flipping to increase dataset diversity and improve model generalization.

Data Mining Techniques:

Association rule mining, clustering, and classification are data mining techniques that can be applied to image captioning systems. Here is a brief overview of how each of these techniques can be used:

- **Association Rule Mining:** Used to identify relationships between different medical terms in reports and corresponding X-ray images. This helps in generating more accurate captions based on frequent co-occurrences.
- **Clustering:** Groups X-ray images based on similarity in features (e.g., disease type, patient age group, or severity). This helps in personalizing image captioning for different patient categories.
- **Classification:** Uses machine learning algorithms to classify X-ray images into different disease categories and generate relevant captions accordingly. For example, a classification model can predict whether an X-ray indicates pneumonia, and the captioning system can describe the findings based on this classification.

Safety and Necessary Precautions:

Ensure that the system follows ethical guidelines and medical standards while handling patient data. Proper validation by medical professionals is essential before deploying the system in real-world clinical settings.

Procedure:

1. Select the domain (Medical Image Processing).
2. Selection of the particular system (Image Captioning for Chest X-Rays).
3. Preprocessing used on the system (image and text processing).
4. Mining techniques applied on the system (association, clustering, classification).

Observation/Program:

In an image captioning system for Chest X-Rays, data mining techniques are used to analyze

large amounts of medical images and generate accurate and meaningful captions, assisting radiologists in diagnosis and reporting.

Conclusion:

Image captioning for Chest X-Rays provides a compelling example of how data mining is an interdisciplinary field, integrating medical expertise, machine learning, and natural language processing to enhance medical diagnostics.

Quiz:

(1) What are the different preprocessing techniques can be applied on dataset?

Preprocessing is a crucial step in preparing a dataset for analysis or machine learning modeling. Here are different preprocessing techniques that can be applied, particularly in the context of **Image Captioning (Chest X-Rays)** or other datasets:

1. Data Cleaning

- Removing duplicate or irrelevant data
- Handling missing values (e.g., imputation, removal)
- Correcting errors in textual data (e.g., spelling correction)
- Eliminating inconsistent or noisy samples

2. Image Preprocessing

- **Resizing:** Standardizing image dimensions to a fixed size for uniform input.
- **Grayscale Conversion:** Converting images to grayscale for better feature extraction.
- **Contrast Enhancement:** Applying histogram equalization or CLAHE to improve visibility.
- **Noise Reduction:** Using filters (e.g., Gaussian, median) to remove noise.
- **Normalization:** Scaling pixel values to a range (0-1) or standardization (mean = 0, std = 1).
- **Data Augmentation:** Rotations, flipping, zooming, and cropping to enhance diversity.

3. Text Preprocessing (For Image Captions)

- **Tokenization:** Splitting sentences into individual words or tokens.
- **Stopword Removal:** Eliminating common but uninformative words (e.g., "the", "is").
- **Lemmatization/Stemming:** Reducing words to their root forms (e.g., "running" → "run").
- **Lowercasing:** Standardizing text case to avoid mismatches.
- **Removing Special Characters & Punctuation:** Ensuring clean and structured data.

4. Feature Extraction

- **CNN-based Feature Extraction:** Using pretrained models (ResNet, VGG, etc.) to extract image features.
- **TF-IDF (Term Frequency-Inverse Document Frequency):** For text-based datasets.
- **Word Embeddings:** Representing words using word vectors (e.g., Word2Vec, GloVe, BERT).

- **Principal Component Analysis (PCA):** Reducing dimensionality while retaining important features.

5. Data Normalization & Standardization

- **Min-Max Scaling:** Rescales values between a defined range (e.g., 0-1).
- **Z-score Normalization:** Transforms data to have mean = 0 and standard deviation = 1.

6. Handling Imbalanced Data

- **Oversampling:** Increasing the count of underrepresented classes (e.g., SMOTE).
- **Undersampling:** Reducing overrepresented class samples to balance the dataset.
- **Class Weighing:** Assigning weights to classes during training.

7. Data Splitting

- Dividing the dataset into training, validation, and testing sets.
- Common splits: **80-10-10** or **70-15-15** for train-validation-test.

(2). What is the use of data mining techniques on particular system?

Data mining enhances **accuracy, efficiency, and automation** in medical image captioning by uncovering patterns in X-ray datasets.

1. Association Rule Mining

Identifies relationships between radiological findings and diseases.

Example: "Lung opacity" frequently appears in pneumonia cases, aiding automated caption generation.

2. Clustering

Groups similar images or patients based on features. Helps personalize captions for different disease categories (e.g., pneumonia, tuberculosis).

3. Classification

Assigns labels to X-rays (e.g., normal, pneumonia, cardiomegaly). CNNs help automate report generation by linking findings to medical conditions.

4. Regression Analysis

Predicts **disease severity** based on past data. Helps generate captions with severity indicators for better diagnosis.

5. Dimensionality Reduction (PCA, t-SNE)

Removes redundant features to improve efficiency. Speeds up training for deep learning models.

6. Anomaly Detection

Identifies rare or misclassified X-rays. Useful for detecting **uncommon diseases** or mislabeled data.

7. NLP for Text Mining

Extracts insights from **radiology reports** to enhance captions. Helps auto-generate text by analyzing common medical phrases.

Suggested References:

1. Han, J., & Kamber, M. (2011). Data mining: concepts and techniques.
2. <https://www.kaggle.com/code/rounakbanik/movie-recommender-systems>

References used by the students:

1. Han, J., & Kamber, M. (2011). Data mining: concepts and techniques.
2. <https://www.kaggle.com/code/ebrahimelgazar/image-captioning-chest-x-rays>

Rubric wise marks obtained:

Rubrics	Knowledge (2)		Problem Recognition (2)		Team Work (2)		Completeness and accuracy (2)		Ethics (2)		Total
	Good (2)	Avg. (1)	Good (2)	Avg. (1)	Good (2)	Avg. (1)	Good (2)	Avg. (1)	Good (2)	Avg. (1)	
Marks											

Experiment No - 2

Aim: Write programs to perform the following tasks of preprocessing (any language).

2.1 Noisy data handling

- Equal Width Binning
- Equal Frequency/Depth Binning

2.2 Normalization Techniques

- Min max normalization
- Z score normalization
- Decimal scaling

2.3. Implement data dispersion measure Five Number Summary generate box plot using python libraries

Date: // Write date of experiment here

Competency and Practical Skills: Programming and statistical methods

Relevant CO: CO1

Objectives:

- (a) To understand Basic Preprocessing Techniques and statistical Measures.
- (b) To show how to implement Preprocessing Techniques.
- (c) To show how to use different Python Libraries to implement Techniques.

Equipment/Instruments: Personal Computer, open-source software for programming

Theory:

2.1 Noisy data handling

Equal Width Binning

Equal Frequency/Depth Binning

Noise: random error or variance in a measured variable

Incorrect attribute values may be due to

- Faulty Data Collection Instruments
- Data Entry Problems
- Data Transmission Problems
- Technology Limitation
- Inconsistency in Naming Convention

Binning: Binning methods smooth a sorted data value by consulting its “neighborhood,” that is, the values around it. The sorted values are distributed into a number of “buckets,” or bins. Because binning methods consult the neighborhood of values, they perform local smoothing.

Data Preprocessing Sorted data for price (in dollars): 4, 8, 15, 21, 21, 24, 25, 28, 34

Partition into (equal-frequency)

bins:

Bin 1: 4, 8, 15
 Bin 2: 21, 21, 24
 Bin 3: 25, 28, 34

Smoothing by bin means:

Bin 1: 9, 9, 9
 Bin 2: 22, 22, 22
 Bin 3: 29, 29, 29

Smoothing by median:

Bin 1: 8, 8, 8
 Bin 2: 21, 21, 21
 Bin 3: 28, 28, 28

Smoothing by bin boundaries:

Bin 1: 4, 4, 15
 Bin 2: 21, 21, 24
 Bin 3: 25, 25, 34

Equal Width Binning :

bins have equal width with a range of each bin are defined as [min + w], [min + 2w] [min + nw]
 where $w = (\max - \min) / (N)$

Example :

5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215

$$W=(215-5)/3=70$$

bin1: 5,10,11,13,15,35,50,55,72 I.e. all values between 5 and 75

bin2: 92 I.e. all values between 75 and 145

bin3: 204,215 I.e. all values between 145 and 215

2.2.Normalization Techniques

Min max normalization

Z score normalization

Decimal scaling

Normalization techniques are used in data preprocessing to scale numerical data to a common range. Here are three commonly used normalization techniques:

The measurement unit used can affect the data analysis. For example, changing measurement units from meters to inches for height, or from kilograms to pounds for weight, may lead to very different results. In general, expressing an attribute in smaller units will lead to a larger range for that attribute, and thus tend to give such an attribute greater effect or “weight.” To help avoid dependence on the choice of measurement units, the data should be normalized or standardized. This involves transforming the data to fall within a smaller or common range such as [-1,1] or [0.0, 1.0]. (The terms standardize and normalize are used interchangeably in data preprocessing, although in statistics, the

latter term also has other connotations.) Normalizing the data attempts to give all attributes an equal weight. Normalization is particularly useful for classification algorithms involving neural networks or distance measurements such as nearest-neighbor classification and clustering. If using the neural network back propagation algorithm for classification mining. There are many methods for data normalization. We Focus on min-max normalization, z-score normalization, and normalization by decimal scaling.

Min-Max Normalization: This technique scales the data to a range of 0 to 1. The formula for min-max normalization is:

$$X_{\text{norm}} = (X - X_{\text{min}}) / (X_{\text{max}} - X_{\text{min}})$$

where X is the original data, X_{min} is the minimum value in the dataset, and X_{max} is the maximum value in the dataset.

Z-Score Normalization: This technique scales the data to have a mean of 0 and a standard deviation of 1. The formula for z-score normalization is:

$$X_{\text{norm}} = (X - X_{\text{mean}}) / X_{\text{std}}$$

where X is the original data, X_{mean} is the mean of the dataset, and X_{std} is the standard deviation of the dataset.

Decimal Scaling: This technique scales the data by moving the decimal point a certain number of places to the left or right. The formula for decimal scaling is:

$$X_{\text{norm}} = X / 10^j$$

where X is the original data and j is the number of decimal places to shift.

2.3. Implement data dispersion measure Five Number Summary generate box plot using python libraries

Five Number Summary

Descriptive Statistics involves understanding the distribution and nature of the data. Five number summary is a part of descriptive statistics and consists of five values and all these values will help us to describe the data.

The minimum value (the lowest value)

25th Percentile or Q1

50th Percentile or Q2 or Median

75th Percentile or Q3

Maximum Value (the highest value)

Let's understand this with the help of an example. Suppose we have some data such as:

11,23,32,26,16,19,30,14,16,10

Here, in the above set of data points our Five Number Summary are as follows:

First of all, we will arrange the data points in ascending order and then calculate the summary: 10,11,14,16,16,19,23,26,30,32

Minimum value: 10

25th Percentile: 14

Calculation of 25th Percentile: $(25/100)*(n+1) = (25/100)*(11) = 2.75$ i.e 3rd value of the data

50th Percentile : 17.5

Calculation of 50th Percentile: $(16+19)/2 = 17.5$

75th Percentile : 26

Calculation of 75th Percentile: $(75/100)*(n+1) = (75/100)*(11) = 8.25$ i.e 8th value of the data

Box plots

Boxplots are the graphical representation of the distribution of the data using Five Number summary values. It is one of the most efficient ways to detect outliers in our dataset.

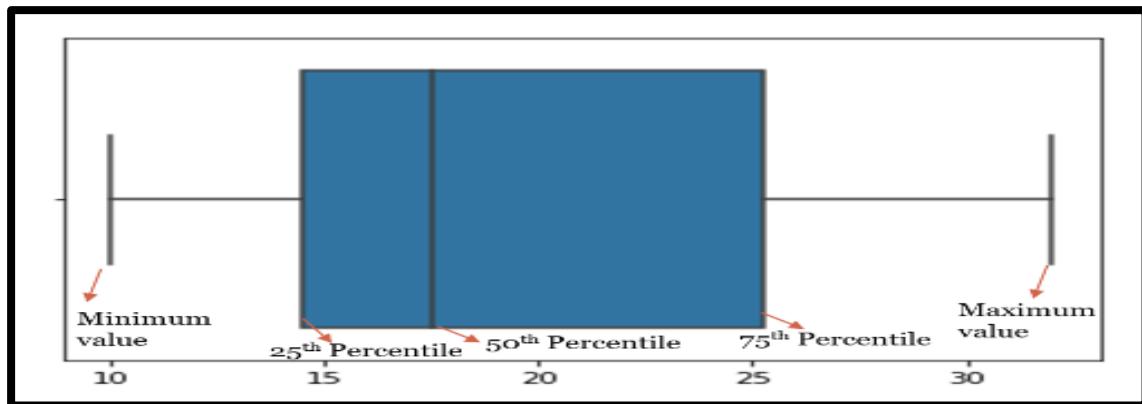


Fig Box plot using Five Number Summary

In statistics, an outlier is a data point that differs significantly from other observations. An outlier may be due to variability in the measurement or it may indicate experimental error; the latter are sometimes excluded from the dataset. An outlier can cause serious problems in statistical analyses.

Safety and necessary Precautions:

Prior to preprocessing, thoroughly clean the data by handling missing values and outliers. Incorrect handling of noisy data can lead to skewed results.

Procedure:

1. Collect raw data from various sourcesSelection of particular system .
2. Handle missing values and outliers.
3. Equal Width Binning:
Divide data into bins of equal width.
4. Equal Frequency/Depth Binning:
Divide data into bins of equal frequency.
5. Min-Max Normalization:
Scale data to a specific range (e.g., [0, 1]).
6. Z-Score Normalization:
Standardize data to have a mean of 0 and standard deviation of 1.
7. Decimal Scaling:
Scale data by a power of 10.
8. Data Dispersion Measure:
Calculate the five-number summary (minimum, Q1, median, Q3, maximum).
9. Box Plot Generation:
Create a box plot to visualize data distribution.

Observation/Program:**Code:**

```
#2.1 Noisy Data Handling
# Generate random Numbers
import random
import statistics
import matplotlib.pyplot as plt
import seaborn as sns

# Generating a larger dataset
random.seed(42)
data = random.sample(range(30, 2000), 30)
data = sorted(data)
print("Random data sample: ", data)

# Number of bins
bins = int(input('Enter the number of bins: '))

# Equal Width Binning
equal_width = []
min_val = min(data)
max_val = max(data)
diff_val = (max_val - min_val) // bins

def range_val(j, limit):
    d = []
    while j < len(data) and data[j] <= limit:
        d.append(data[j])
        j += 1
    return j, d

j = 0
for i in range(1, bins+1):
    j, val = range_val(j, min_val + (i * diff_val))
    equal_width.append(val)

print("Equal Width : ")
print(equal_width)

# Equal Frequency Binning
equal_freq = []
size_bin = len(data) // bins

for i in range(1, bins+1):
    start = (i-1) * size_bin
    stop = start + size_bin
    equal_freq.append(data[start: stop])

print("Equal Frequency : ")
print(equal_freq)

# Smoothing Techniques
```

```

def smooth_mean(data):
    return [[statistics.mean(b)] * len(b) for b in data]

def smooth_median(data):
    return [[statistics.median(b)] * len(b) for b in data]

def smooth_bound(data):
    return [[b[0]] + [min(b) if (x - min(b)) <= (max(b) - x) else max(b) for x in b[1:-1]] + [b[-1]]]
for b in data]

print("Smooth mean for Equal Frequency:", smooth_mean(equal_freq))
print("Smooth mean for Equal Width:", smooth_mean(equal_width))
print("Smooth median for Equal Frequency:", smooth_median(equal_freq))
print("Smooth median for Equal Width:", smooth_median(equal_width))
print("Smooth bound for Equal Frequency:", smooth_bound(equal_freq))
print("Smooth bound for Equal Width:", smooth_bound(equal_width))

```

2.2 Normalization Techniques

- 1) Min max normalization
- 2) Z score normalization
- 3) Decimal scaling

Code:

```

# 2.2 Normalization Techniques
# Min max normalization
# Z score normalization
# Decimal scaling

data_min = min(data)
data_max = max(data)
data_mean = statistics.mean(data)
data_std = statistics.stdev(data)
abs_max = max(abs(min(data)), abs(max(data)))
decimal_places = len(str(abs_max))

# 1. min-max normalization

def min_max(val, new_min=0.0, new_max=1.0):
    return round(((val - data_min) / (data_max - data_min)) * (new_max - new_min) + new_min, 2)

min_max_norm = [min_max(i) for i in data]
print("Min-Max Normalization:", min_max_norm)

# 2. Z score normalization

def z_score(val):
    return round((val - data_mean) / data_std, 2)

```

```
z_norm = [z_score(i) for i in data]
print("Z-score Normalization:", z_norm)
```

3. Decimal Scaling

```
def dec_scale(val):
    return round(val / (10 ** decimal_places), 2)

dec_norm = [dec_scale(i) for i in data]
print("Decimal Scaling Normalization:", dec_norm)
```

2.3 Implement data dispersion measure Five Number Summary generate box plot using python libraries.

```
Q1 = statistics.median(data[:len(data)//2])
Q2 = statistics.median(data)
Q3 = statistics.median(data[len(data)//2:])
```

```
# Printing Five Number Summary
print("Five Number Summary")
print("Minimum:", data_min)
print("Q1 (25%):", Q1)
print("Q2 (50%) (Median):", Q2)
print("Q3 (75%):", Q3)
print("Maximum:", data_max)
```

```
# Box Plot Visualization (Now with values along x-axis)
plt.figure(figsize=(8, 5))
sns.boxplot(x=data) # Changed from vertical to horizontal box plot
plt.title("Box Plot of the Dataset (X-Axis)")
plt.xlabel("Values")
plt.show()
```

2.1 Noisy Data Handling

Output:

2.2 Normalization Techniques

Output:

```

Min-Max Normalization: [0.0, 0.0, 0.01, 0.01, 0.07, 0.08, 0.09, 0.1, 0.13, 0.2, 0.22, 0.23, 0.24, 0.25, 0.29, 0.46, 0.55, 0.6, 0.62, 0.65, 0.66, 0.71, 0.72, 0.75, 0.78, 0.8, 0.82, 0.82, 0.83, 1.0]
Z-score Normalization: [-1.31, -1.31, -1.3, -1.29, -1.09, -1.07, -1.04, -1.0, -0.9, -0.69, -0.62, -0.6, -0.57, -0.53, -0.42, 0.11, 0.4, 0.55, 0.61, 0.71, 0.75, 0.89, 0.92, 1.02, 1.11, 1.16, 1.23, 1.25, 1.25, 1.79]
Decimal Scaling Normalization: [0.01, 0.01, 0.01, 0.01, 0.02, 0.02, 0.02, 0.03, 0.03, 0.04, 0.05, 0.05, 0.05, 0.05, 0.06, 0.09, 0.11, 0.11, 0.12, 0.12, 0.13, 0.13, 0.14, 0.14, 0.15, 0.15, 0.15, 0.15, 0.19]

```

2.3 Implement data dispersion measure Five Number Summary

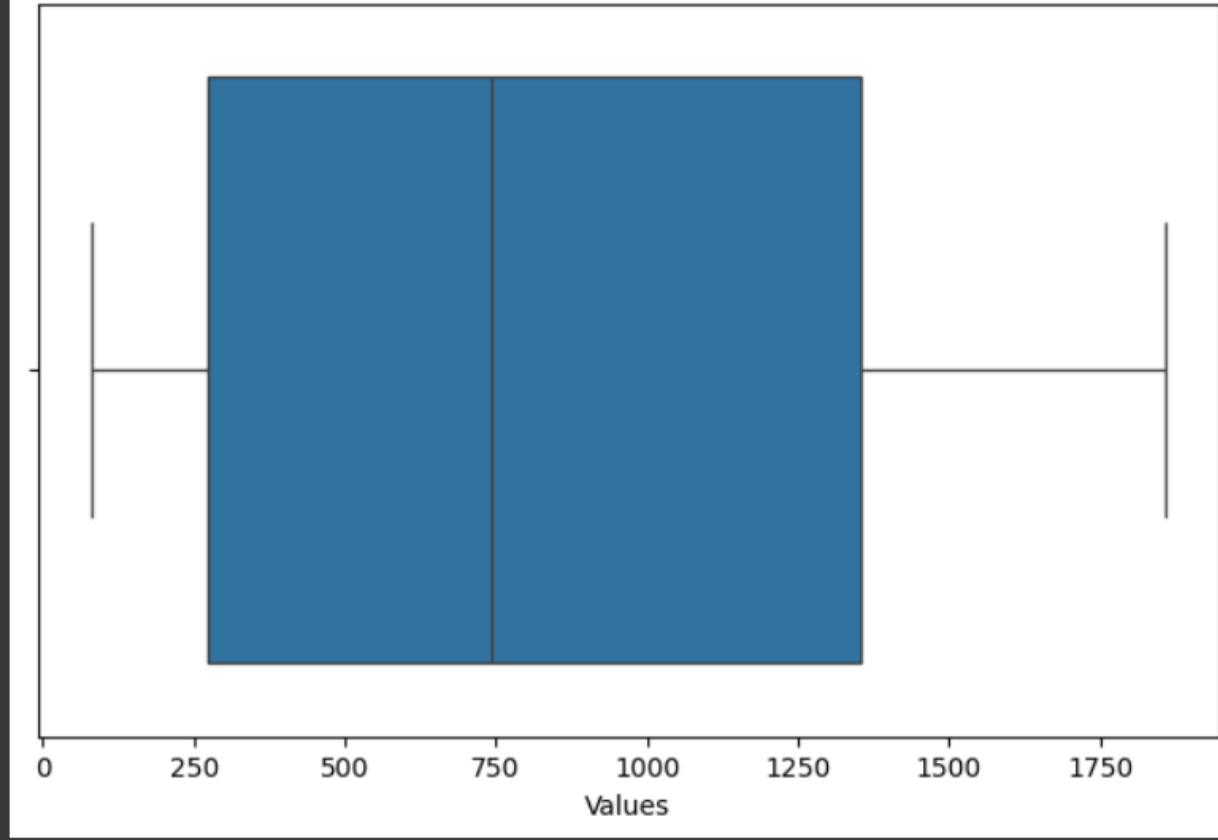
Output:

```

Five Number Summary
Minimum: 81
Q1 (25%): 258
Q2 (50%) (Median): 743.5
Q3 (75%): 1360
Maximum: 1857

```

Box Plot of the Dataset (X-Axis)



Conclusion:

Binning, Normalization techniques and the five number summary are both important tools in data preprocessing that help prepare data for data mining tasks.

Quiz:

(1) What is Five Number summary? How to generate box plot using Python Libraries?

A **Five Number Summary** is a statistical measure that provides a concise description of a dataset's distribution. It consists of five key values:

1. **Minimum** – The smallest value in the dataset.
2. **Q1 (First Quartile / 25th Percentile)** – The median of the lower half of the data.
3. **Q2 (Median / 50th Percentile)** – The middle value of the dataset.

4. **Q3 (Third Quartile / 75th Percentile)** – The median of the upper half of the data.
5. **Maximum** – The largest value in the dataset.

The Five Number Summary helps to understand the spread and central tendency of the data. It is also useful for detecting outliers and skewness in the dataset.

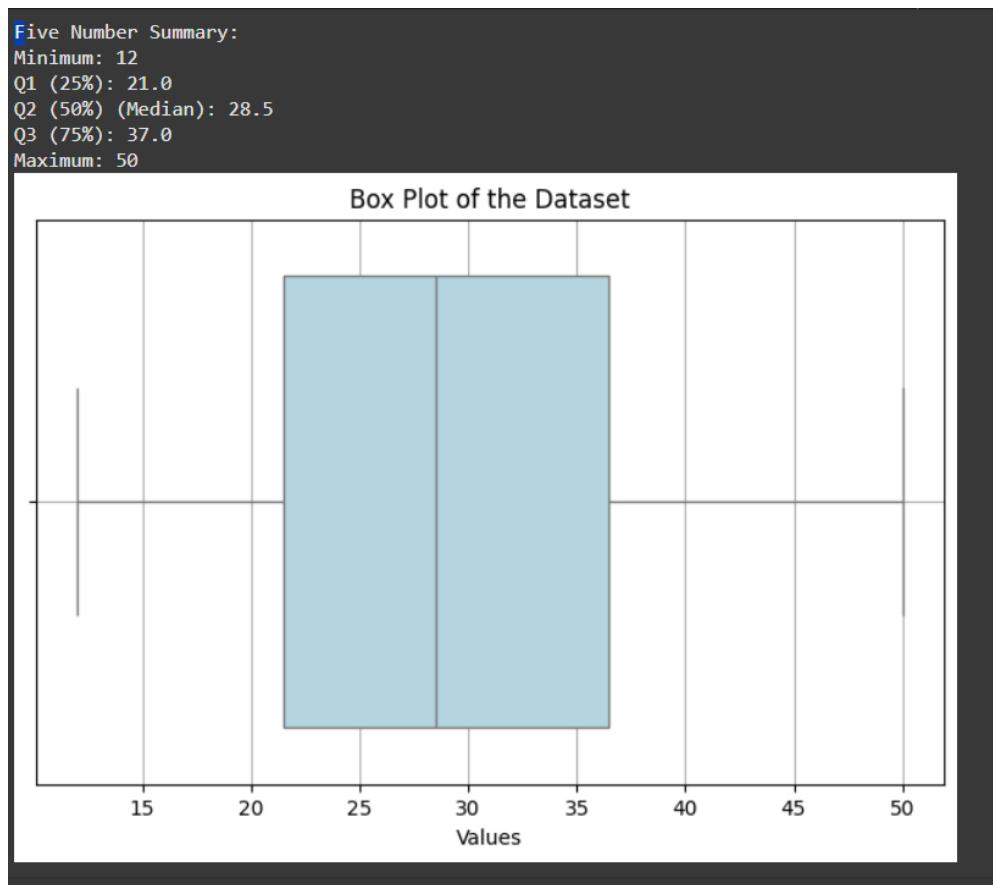
```
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import statistics

# Sample dataset
data = [12, 25, 36, 15, 40, 50, 22, 30, 18, 27, 35, 45, 20, 28, 32, 38, 48, 16, 26, 29]
data.sort()

# Calculating Five Number Summary
data_min = min(data)
Q1 = statistics.median(data[:len(data)//2]) # 25th percentile
Q2 = statistics.median(data) # Median (50th percentile)
Q3 = statistics.median(data[len(data)//2:]) # 75th percentile
data_max = max(data)

# Printing Five Number Summary
print("Five Number Summary:")
print("Minimum:", data_min)
print("Q1 (25%):", Q1)
print("Q2 (50%) (Median):", Q2)
print("Q3 (75%):", Q3)
print("Maximum:", data_max)

# Generating Box Plot
plt.figure(figsize=(8, 5))
sns.boxplot(x=data, color="lightblue") # Box plot with values along x-axis
plt.title("Box Plot of the Dataset")
plt.xlabel("Values")
plt.grid(True)
plt.show()
```



(2) What is Normalization techniques?

Normalization is a data preprocessing technique used to scale numerical data into a specific range, typically **[0,1]** or **[-1,1]**. It ensures that no particular feature dominates the model due to larger numerical values. Normalization is especially useful in **machine learning, data mining, and statistical analysis** when using algorithms that rely on distance calculations (e.g., k-NN, clustering, and neural networks).

Types of Normalization:

1. Min-Max Normalization

- Scales data between **0 and 1** using:

$$X_{\text{norm}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

2. Z-Score Normalization (Standardization)

Transforms data to have **mean = 0** and **standard deviation = 1** using:

$$X_{\text{norm}} = \frac{X - X_{\text{mean}}}{X_{\text{std}}}$$

3. Decimal Scaling Normalization

Scales data by dividing by 10^j , where j is the smallest integer ensuring all values fall in [-1,1]:

$$X_{\text{norm}} = \frac{X}{10^j}$$

(3) What are the different smoothing techniques?

Smoothing techniques are used in **data preprocessing** to reduce **noise** and **variability** in a dataset while retaining its essential patterns. These techniques help improve data quality, making it easier for **statistical analysis, machine learning, and visualization**.

Technique	Best For	Key Feature
Binning	Numerical data	Groups data into bins
Moving Average	Time-series	Averages over a window
Exponential Smoothing	Forecasting	Weights decrease exponentially
Gaussian Smoothing	Image processing	Uses Gaussian function
Regression Smoothing	Trend detection	Fits curves to data

Suggested Reference:

J. Han, M. Kamber, "Data Mining Concepts and Techniques", Morgan Kaufmann

References used by the students:

// Write references used by you here

Rubric wise marks obtained:

Rubrics	Knowledge (2)		Problem Recognition (2)		Logic Building (2)		Completeness and accuracy (2)		Ethics (2)		Total
	Good (2)	Avg. (1)	Good (2)	Avg. (1)	Good (2)	Avg. (1)	Good (2)	Avg. (1)	Good (2)	Avg. (1)	
Marks											

Experiment No - 3

Aim: To perform hand on experiments of data preprocessing with sample data on Orange tool.

Date: // Write date of experiment here

Competency and Practical Skills: Exploration and Understanding of Tool

Relevant CO: CO1 & CO4

Objectives: 1) improve users' understanding of data preprocessing techniques
2) Familiarize them with the tool

Equipment/Instruments: Orange tool

Safety and necessary Precautions:

Document the steps you take in Orange, including the specific preprocessing techniques, parameters used, and the reasoning behind your choices.

Procedure:

1. Install and set up Orange.
2. Import your dataset.
3. Apply different pre processing methods .

Demonstration of Tool:

Data Preprocessing With Orange tool

Preprocesses data with selected methods. Inputs Data: input dataset Outputs Preprocessor: preprocessing method Preprocessed Data: data preprocessed with selected methods Preprocessing is crucial for achieving better-quality analysis results. The Preprocess widget offers several preprocessing methods that can be combined in a single preprocessing pipeline. Some methods are available as separate widgets, which offer advanced techniques and greater parameter tuning.

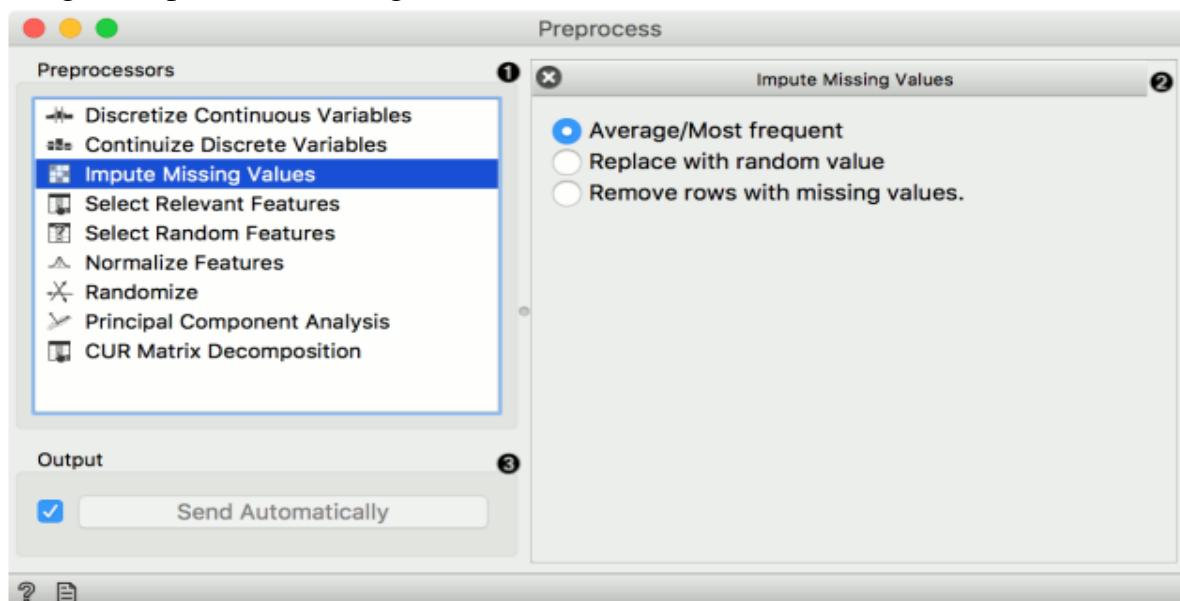
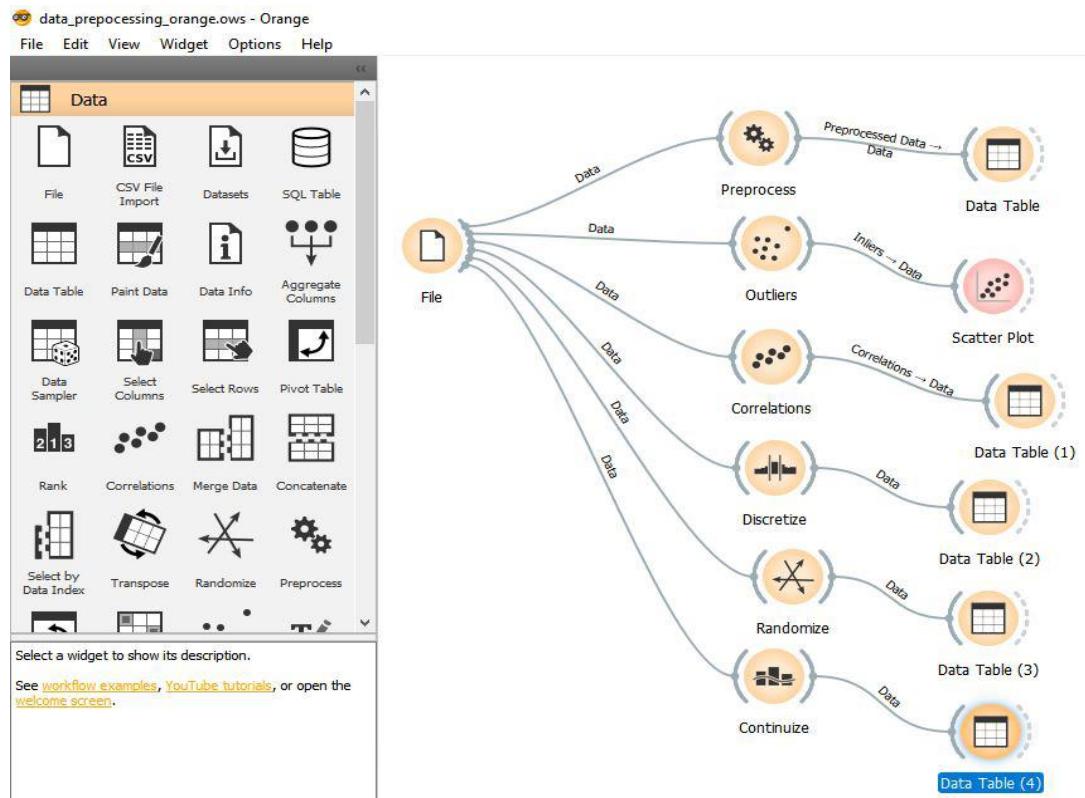
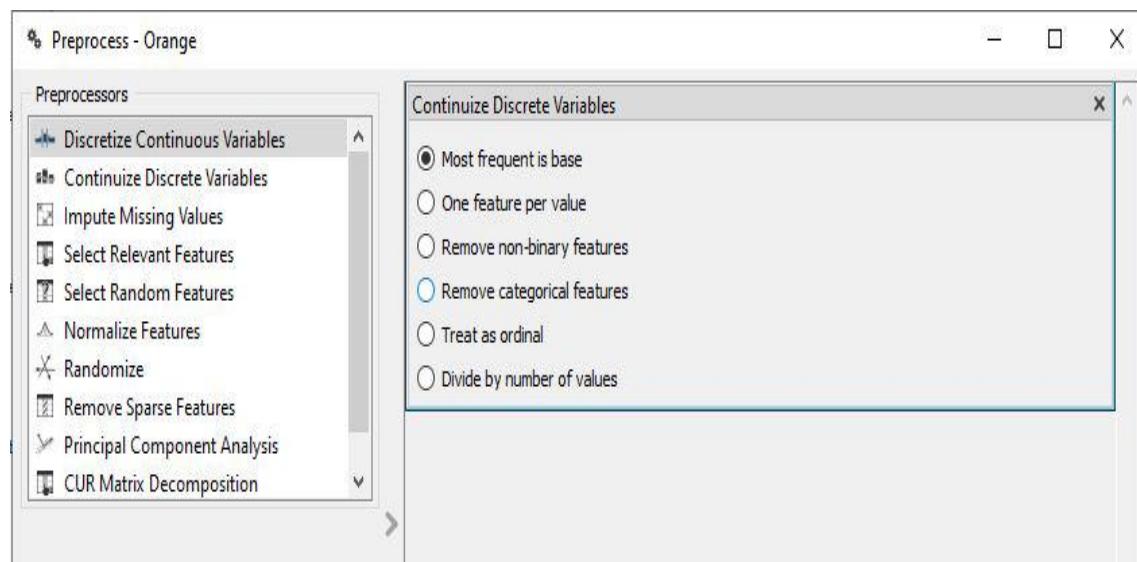


Fig Handling Missing Values

1. List of preprocessors. Double click the preprocessors you wish to use and shuffle their order by dragging them up or down. You can also add preprocessors by dragging them from the left menu to the right.
2. Preprocessing pipeline.
3. When the box is ticked (Send Automatically), the widget will communicate changes automatically. Alternatively, click Send.



➤ Preprocessed Technique :



[Fig: Descrete Continuous variables -> Most Frequent is base Used]

➤ **Data Table of Preprocessed Data**

Data Table - Orange

Info
150 instances (no missing data)
4 features
Target with 3 values
No meta attributes

Variables
 Show variable labels (if present)
 Visualize numeric values
 Color by instance classes

Selection
 Select full rows

Restore Original Order
 Send Automatically

	iris	sepal length	sepal width	petal length	petal width
1	Iris-setosa	5.1	3.5	1.4	0.2
2	Iris-setosa	4.9	3.0	1.4	0.2
3	Iris-setosa	4.7	3.2	1.3	0.2
4	Iris-setosa	4.6	3.1	1.5	0.2
5	Iris-setosa	5.0	3.6	1.4	0.2
6	Iris-setosa	5.4	3.9	1.7	0.4
7	Iris-setosa	4.6	3.4	1.4	0.3
8	Iris-setosa	5.0	3.4	1.5	0.2
9	Iris-setosa	4.4	2.9	1.4	0.2
10	Iris-setosa	4.9	3.1	1.5	0.1
11	Iris-setosa	5.4	3.7	1.5	0.2
12	Iris-setosa	4.8	3.4	1.6	0.2
13	Iris-setosa	4.8	3.0	1.4	0.1
14	Iris-setosa	4.3	3.0	1.1	0.1
15	Iris-setosa	5.8	4.0	1.2	0.2
16	Iris-setosa	5.7	4.4	1.5	0.4
17	Iris-setosa	5.4	3.9	1.3	0.4
18	Iris-setosa	5.1	3.5	1.4	0.3
19	Iris-setosa	5.7	3.8	1.7	0.3
20	Iris-setosa	5.1	3.8	1.5	0.3
21	Iris-setosa	5.4	3.4	1.7	0.2
22	Iris-setosa	5.1	3.7	1.5	0.4
23	Iris-setosa	4.6	3.6	1.0	0.2
24	Iris-setosa	5.1	3.3	1.7	0.5
25	Iris-setosa	4.8	3.4	1.9	0.2
26	Iris-setosa	5.0	3.0	1.6	0.2
27	Iris-setosa	5.0	3.4	1.6	0.4
28	Iris-setosa	5.2	3.5	1.5	0.2
29	Iris-setosa	5.2	3.4	1.4	0.2
30	Iris-setosa	4.7	3.2	1.6	0.2
31	Iris-setosa	4.8	3.1	1.6	0.2
32	Iris-setosa	5.4	3.4	1.5	0.4
33	Iris-setosa	5.2	4.1	1.5	0.1
34	Iris-setosa	5.5	4.2	1.4	0.2

Observation:

1. Load the Dataset

- **Tool Used:** CSV File Import
- **Dataset:** Titanic dataset
- **Reasoning:** Import the dataset into Orange for preprocessing.

2. Remove Unnecessary Columns

- **Tool Used:** Select Columns
- **Columns Removed:**
 - PassengerId
 - Name
 - Ticket
 - Cabin
- **Reasoning:** These columns are irrelevant to the survival prediction and do not contribute useful information.

3. Handle Missing Values

- **Tool Used:** Impute
- **Imputation Strategy:**
 - Age → Mean
 - Embarked → Mean
- **Reasoning:** Filling missing values with the mean ensures numerical stability without dropping rows.

4. Data Visualization

- **Tool Used:** Distributions
- **Visualization of:** Gender-based survival rates
- **Reasoning:** To analyze how survival varies across different groups.

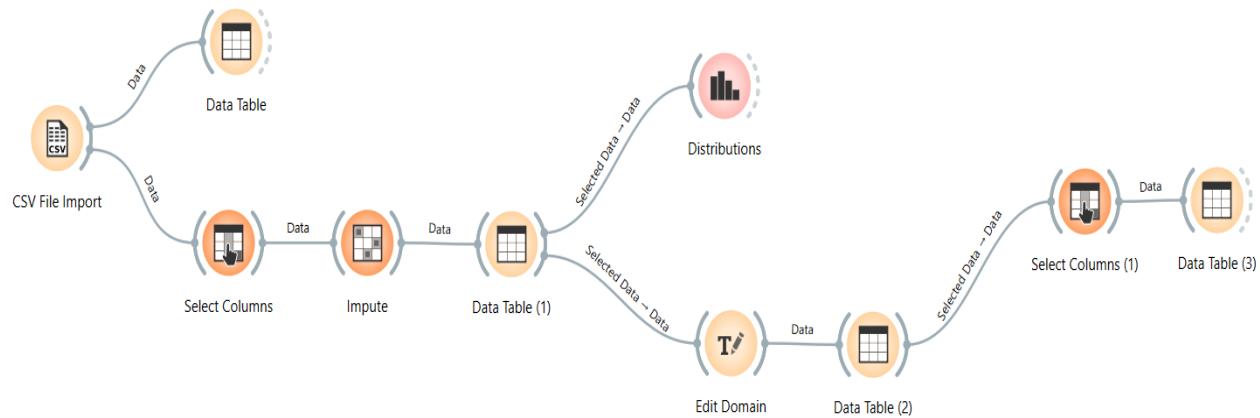
5. Encode Categorical Variables

- **Tool Used:** Edit Domain
- **Changes Made:**
 - Sex:
 - Male → 0
 - Female → 1
 - Embarked:
 - S → 0
 - C → 1
 - Q → 2
- **Reasoning:** Converting categorical variables into numerical values makes them usable for machine learning models.

6. Separate Target Variable

- Tool Used:** Select Columns
- Target Variable:** Survived
- Feature Variables:** All other relevant columns
- Reasoning:** Separating the target variable allows proper supervised learning processing.

Output:



Data before preprocessing:

Info
891 instances
9 features (2.2 % missing data)
No target variable.
3 meta attributes (25.7 % missing data)

Variables
 Show variable labels (if present)
 Visualize numeric values
 Color by instance classes

Selection
 Select full rows

Restore Original Order
 Send Automatically

	Name	Ticket	Cabin	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
1	Braund, Mr. Owl...	A/5 21171	?	1	0	3	male	22	1	0	7.25	S
2	Cumings, Mrs. J...	PC 17599	C85	2	1	1	female	38	1	0	71.2833	C
3	Heikkinen, Miss...	STON/O2. 3101...	?	3	1	3	female	26	0	0	7.925	S
4	Futrelle, Mrs. Ja...	113803	C123	4	1	1	female	35	1	0	53.1	S
5	Allen, Mr. Willia...	373450	?	5	0	3	male	35	0	0	8.05	S
6	Moran, Mr. Jam...	330877	?	6	0	3	male	?	0	0	8.4583	Q
7	McCarthy, Mr. T...	17463	E46	7	0	1	male	54	0	0	51.8625	S
8	Palsson, Master....	349909	?	8	0	3	male	2	3	1	21.075	S
9	Johnson, Mrs.	347742	?	9	1	3	female	27	0	2	11.1333	S
10	Nasser, Mrs. Ni...	237736	?	10	1	2	female	14	1	0	30.0708	C
11	Sandstrom, Mis...	PP 9549	G6	11	1	3	female	4	1	1	16.7	S
12	Bonnell, Miss. El...	113783	C103	12	1	1	female	58	0	0	26.55	S
13	Saudercock, ...	A/5. 2151	?	13	0	3	male	20	0	0	8.05	S
14	Andersson, Mr. ...	347082	?	14	0	3	male	39	1	5	31.275	S
15	Vestrom, Miss. ...	350406	?	15	0	3	female	14	0	0	7.8542	S
16	Hewlett, Mrs. (...)	248706	?	16	1	2	female	55	0	0	16	S
17	Rice, Master. Eu...	382652	?	17	0	3	male	2	4	1	29.125	Q
18	Williams, Mr. C...	244373	?	18	1	2	male	?	0	0	13	S
19	Vander Planke, ...	345763	?	19	0	3	female	31	1	0	18	S
20	Masselmani, Mr. ...	2649	?	20	1	3	female	?	0	0	7.225	C
21	Fynney, Mr. Jos...	239865	?	21	0	2	male	35	0	0	26	S
22	Beesley, Mr. La...	248698	D56	22	1	2	male	34	0	0	13	S
23	McGowan, Miss. ...	330923	?	23	1	3	female	15	0	0	8.0292	Q
24	Sloper, Mr. Willi...	113788	A6	24	1	1	male	28	0	0	35.5	S
25	Palsson, Miss. T...	349909	?	25	0	3	female	8	3	1	21.075	S
26	Asplund, Mrs. C...	347077	?	26	1	3	female	38	1	5	31.3875	S
27	Emir, Mr. Farred...	2631	?	27	0	3	male	?	0	0	7.225	C
28	Fortune, Mr. Ch...	19950	C23 C25 C27	28	0	1	male	19	3	2	263	S
29	O'Dwyer, Miss. ...	330959	?	29	1	3	female	?	0	0	7.8792	Q
30	Todoroff, Mr. La...	349216	?	30	0	3	male	?	0	0	7.8958	S
	Hauck, Mr. Dan...	3417601	?	31	0	1	male	40	0	0	27.7208	C

DataTable of Preprocessed Data:



Conclusion: Orange is a powerful open-source data analysis and visualization tool for machine learning and data mining tasks. It provides a wide variety of functionalities including data visualization, data preprocessing, feature selection, classification, regression, clustering, and more. Its user-friendly interface and drag-and-drop workflow make it easy for non-experts to work with and understand machine learning concepts.

Quiz:

(1) What is the purpose of Orange's Preprocess method?

The **Preprocess** method in **Orange** is used to clean, transform, and prepare data for analysis and machine learning. It helps with tasks like:

- Handling missing values (imputation)
- Scaling and normalization
- Encoding categorical variables
- Feature selection and extraction

(2) What is the use of orange tool?

Orange is a **data visualization and machine learning tool** used for:

- Data preprocessing (cleaning, transformation)
- Exploratory data analysis (visualization, statistics)
- Machine learning (classification, clustering, regression)
- Text and network analysis
- It provides an easy drag-and-drop interface for both beginners and advanced users to analyze data without coding.

Suggested Reference:

1. J. Han, M. Kamber, "Data Mining Concepts and Techniques", Morgan Kaufman
2. <https://orangedatamining.com/docs/>

References used by the students:

// Write references used by you here

Rubric wise marks obtained:

Rubrics	Knowledge (2)		Problem Recognition (2)		Tool Usage/ Demonstrati on (2)		Communicati on Skill (2)		Ethics (2)		Total
	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	
Marks											

Experiment No - 4

Aim: Implement Apriori algorithm of association rule data mining technique in any Programming language.

Date:

Competency and Practical Skills: Logic building, Programming and Analyzing

Relevant CO: CO2

Objectives: To implement basic logic for association rule mining algorithm with support and confidence measures.

Equipment/Instruments: Personal Computer, open-source software for programming

Theory:

The Apriori algorithm is a classic and fundamental data mining algorithm used for discovering association rules in transactional datasets.

- Apriori is designed for finding associations or relationships between items in a dataset. It's commonly used in market basket analysis and recommendation systems.
- Apriori discovers frequent itemsets, which are sets of items that frequently co-occur in transactions. A frequent itemset is a set of items that appears in a minimum number of transactions, known as the "support threshold."
- Support and Confidence: Support measures how often an itemset appears in the dataset, while confidence measures how often a rule is true. High-confidence rules derived from frequent itemsets are of interest.

$$\text{Support (A)} = \frac{\text{Number of transaction in which A appears}}{\text{Total number of transactions}}$$

$$\text{Confidence (A} \rightarrow \text{B)} = \frac{\text{Support(AUB)}}{\text{Support(A)}}$$

- Apriori uses an iterative approach to progressively discover frequent itemsets of increasing size. It starts with finding frequent 1-itemsets, then 2-itemsets, and so on.
- The algorithm employs pruning techniques to reduce the number of candidate itemsets that need to be checked, making it more efficient.
- Apriori is widely used in retail for market basket analysis. It helps retailers understand which products are often purchased together, allowing for optimized store layouts, targeted marketing, and product recommendations.

Safety and necessary Precautions:

Ensure that your dataset is clean and free from missing values, outliers, and inconsistencies.

1. **Procedure:**
2. Import the dataset that you want to analyze for association rules.
3. Define the minimum support and confidence thresholds for the Apriori algorithm. These parameters control the minimum occurrence of itemsets and the minimum confidence level for rules.
4. Implement the Apriori algorithm to discover frequent itemsets.
5. Use the frequent itemsets obtained from the previous step to generate association rules

Observation/Program:

```

import pandas as pd
from itertools import combinations

# Load dataset
file_path = "basket_analysis.csv"
df = pd.read_csv(file_path).drop(columns=["Unnamed: 0"]) # Remove index column

# Convert dataset into list of transactions
transactions = [set(df.index[row]) for _, row in df.iterrows()]

# Function to calculate support
def support(itemset):
    return sum(1 for t in transactions if itemset.issubset(t)) / len(transactions)

# Minimum support and confidence thresholds
min_support, min_confidence = 0.05, 0.5

# Generate frequent itemsets
frequent_itemsets = {frozenset([item]): support(frozenset([item])) for item in df.columns}
frequent_itemsets = {k: v for k, v in frequent_itemsets.items() if v >= min_support}

k = 2
while True:
    candidates = {a | b for a in frequent_itemsets for b in frequent_itemsets if len(a | b) == k}
    candidate_supports = {c: support(c) for c in candidates if support(c) >= min_support}
    if not candidate_supports:
        break
    frequent_itemsets.update(candidate_supports)
    k += 1

# Generate association rules
rules = []
for itemset in frequent_itemsets:

```

```

for i in range(1, len(itemset)):
    for antecedent in map(frozenset, combinations(itemset, i)):
        consequent = itemset - antecedent
        conf = support(itemset) / support(antecedent) if support(antecedent) > 0 else 0
        if conf >= min_confidence:
            rules.append((antecedent, consequent, conf))

# Display top 10 rules
for ant, con, conf in sorted(rules, key=lambda x: x[2], reverse=True)[:10]:
    print(f"{{set(ant)}} -> {{set(con)}} (Confidence: {conf:.2f})")

```

Observations:

```

{'Milk', 'Dill', 'Unicorn'} -> {'chocolate'} (Confidence: 0.68)
{'Cheese', 'Sugar', 'Unicorn'} -> {'Kidney Beans'} (Confidence: 0.67)
{'Cheese', 'Ice cream', 'Yogurt'} -> {'Kidney Beans'} (Confidence: 0.66)
{'Milk', 'Cheese', 'Dill'} -> {'chocolate'} (Confidence: 0.65)
{'Apple', 'Corn', 'Onion'} -> {'Sugar'} (Confidence: 0.65)
{'Milk', 'Nutmeg', 'Corn'} -> {'Kidney Beans'} (Confidence: 0.65)
{'Kidney Beans', 'Dill', 'Onion'} -> {'Cheese'} (Confidence: 0.65)
{'Cheese', 'Dill', 'Unicorn'} -> {'chocolate'} (Confidence: 0.65)
{'Butter', 'Dill', 'Unicorn'} -> {'chocolate'} (Confidence: 0.65)
{'Dill', 'Unicorn', 'Onion'} -> {'chocolate'} (Confidence: 0.65)

```

Conclusion:

The implementation of the Apriori algorithm successfully identified frequent itemsets and strong association rules in the dataset. The results highlight key product relationships based on transactional data. For instance, items like chocolate, butter, and yogurt appeared frequently in transactions, indicating high customer demand. Additionally, association rules such as {Dill, Milk, Unicorn} → {Chocolate} with a high confidence of 68.1% suggest strong co-purchasing trends.

Quiz:

(1) What Do you Mean by Association rule mining?

Association Rule Mining is a data mining technique used to discover interesting relationships or patterns (associations) between items in large datasets. It identifies rules of the form "**If X, then Y**", where X and Y are itemsets, based on support and confidence measures. It is widely used in market basket analysis, recommendation systems, and business analytics.

(2) What are the different measures are used in apriori algorithm?

The Apriori algorithm uses the following key measures:

Support – The frequency of an itemset appearing in transactions.

Support(X)=Transactions containing X/Total transactions

Confidence – The likelihood that item Y is bought when item X is bought.
 $\text{Confidence}(X \rightarrow Y) = \text{Support}(X \cup Y) / \text{Support}(X)$

Lift – Measures how much more likely Y is purchased when X is present, compared to random chance.
 $\text{Lift}(X \rightarrow Y) = \text{Confidence}(X \rightarrow Y) / \text{Support}(Y)$

Suggested Reference:

- J. Han, M. Kamber, “Data Mining Concepts and Techniques”, Morgan Kaufmann

References used by the students:

<https://www.geeksforgeeks.org>

Rubric wise marks obtained:

Rubrics	Knowledge (2)		Problem Recognition (2)		Logic Building (2)		Completeness and accuracy (2)		Ethics (2)		Total
	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	
Marks											

Experiment No - 5

Aim: Apply association rule data mining technique on sample data sets using WEKA.

Date:

Competency and Practical Skills: Exploration and Understanding of Tool

Relevant CO: CO2 & CO4

Objectives: 1) improve users' understanding of Association rule mining techniques
2) Familiarize with the tool

Equipment/Instruments: WEKA Tool

Safety and necessary Precautions:

Properly evaluate the generated association rules and avoid drawing incorrect conclusions.
Consider using various rule quality metrics to assess rule significance.

Procedure:

1. Install and set up WEKA.
2. Import your dataset.
3. Apply association rule data mining methods .

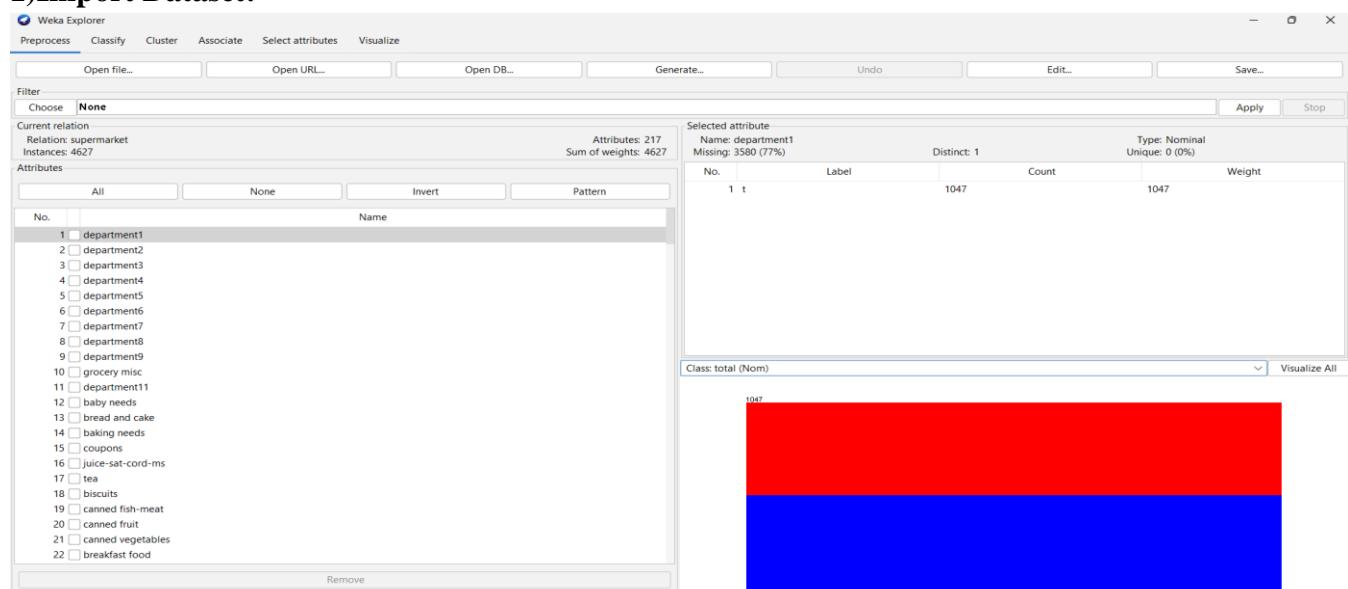
Demonstration of Tool:

The Apriori Algorithm in WEKA identifies frequent itemsets and generates association rules from the dataset.

It uses support and confidence measures to determine rule significance.

Observations:

1) Import Dataset:



2)Apriori Method:

The screenshot shows the Weka Explorer interface with the 'Associate' tab selected. Under 'Choose', 'Apriori' is selected with parameters: -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S 1.0 -c -1. The 'Start' button is highlighted. The 'Result list' pane shows the start time '21:48:59 - Apriori'. The 'Associator output' pane displays the Apriori process details, including minimum support (0.15), minimum confidence (0.9), and number of cycles (17). It lists generated itemsets L(1) through L(6) and their sizes. The 'Best rules found:' section contains 10 rules, each with items, confidence, lift, and support values.

```

Apriori
=====
Minimum support: 0.15 (694 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 17

Generated sets of large itemsets:

Size of set of large itemsets L(1): 44
Size of set of large itemsets L(2): 380
Size of set of large itemsets L(3): 910
Size of set of large itemsets L(4): 633
Size of set of large itemsets L(5): 105
Size of set of large itemsets L(6): 1

Best rules found:

1. biscuits=t frozen foods=t total=high 788 ==> bread and cake=t 723 <conf:(0.92)> lift:(1.27) lev:(0.03) [155] conv:(3.35)
2. baking needs=t biscuits=t fruit=t total=high 760 ==> bread and cake=t 696 <conf:(0.92)> lift:(1.27) lev:(0.03) [149] conv:(3.28)
3. baking needs=t frozen foods=t fruit=t total=high 770 ==> bread and cake=t 705 <conf:(0.92)> lift:(1.27) lev:(0.03) [150] conv:(3.27)
4. biscuits=t fruit=t vegetables=t total=high 815 ==> bread and cake=t 746 <conf:(0.92)> lift:(1.27) lev:(0.03) [159] conv:(3.26)
5. party snack foods=t fruit=t total=high 854 ==> bread and cake=t 779 <conf:(0.91)> lift:(1.27) lev:(0.04) [164] conv:(3.15)
6. biscuits=t frozen foods=t vegetables=t total=high 797 ==> bread and cake=t 725 <conf:(0.91)> lift:(1.26) lev:(0.03) [151] conv:(3.06)
7. baking needs=t biscuits=t vegetables=t total=high 772 ==> bread and cake=t 701 <conf:(0.91)> lift:(1.26) lev:(0.03) [145] conv:(3.01)
8. biscuits=t fruit=t total=high 954 ==> bread and cake=t 866 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(3)
9. frozen foods=t fruit=t vegetables=t total=high 834 ==> bread and cake=t 757 <conf:(0.91)> lift:(1.26) lev:(0.03) [156] conv:(3)
10. frozen foods=t fruit=t total=high 969 ==> bread and cake=t 877 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(2.92)

```

3)FPGrowth:

The screenshot shows the Weka Explorer interface with the 'Associate' tab selected. Under 'Choose', 'FPGrowth' is selected with parameters: -P 2 -I -1 -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1. The 'Start' button is highlighted. The 'Result list' pane shows the start times '21:48:59 - Apriori' and '21:51:46 - FPGrowth'. The 'Associator output' pane displays run information, including scheme (weka.associations.FPGrowth), relation (supermarket), instances (4627), and attributes (217). It also shows the 'Associator model (full training set)' and lists 16 rules found by FPGrowth, which are identical to the Apriori rules above.

```

Scheme: weka.associations.FPGrowth -P 2 -I -1 -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1
Relation: supermarket
Instances: 4627
Attributes: 217
[list of attributes omitted]
Associator model (full training set)

FPGrowth found 16 rules (displaying top 10)

1. [fruit=t, frozen foods=t, biscuits=t, total=high]: 788 ==> [bread and cake=t]: 723 <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.35)
2. [fruit=t, baking needs=t, biscuits=t, total=high]: 760 ==> [bread and cake=t]: 696 <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.28)
3. [fruit=t, baking needs=t, frozen foods=t, total=high]: 770 ==> [bread and cake=t]: 705 <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.27)
4. [fruit=t, vegetables=t, biscuits=t, total=high]: 815 ==> [bread and cake=t]: 746 <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.26)
5. [fruit=t, party snack foods=t, total=high]: 854 ==> [bread and cake=t]: 779 <conf:(0.91)> lift:(1.27) lev:(0.04) conv:(3.15)
6. [vegetables=t, frozen foods=t, biscuits=t, total=high]: 797 ==> [bread and cake=t]: 725 <conf:(0.91)> lift:(1.26) lev:(0.03) conv:(3.06)
7. [vegetables=t, baking needs=t, biscuits=t, total=high]: 772 ==> [bread and cake=t]: 701 <conf:(0.91)> lift:(1.26) lev:(0.03) conv:(3.01)
8. [fruit=t, biscuits=t, total=high]: 954 ==> [bread and cake=t]: 866 <conf:(0.91)> lift:(1.26) lev:(0.04) conv:(3)
9. [fruit=t, vegetables=t, frozen foods=t, total=high]: 834 ==> [bread and cake=t]: 757 <conf:(0.91)> lift:(1.26) lev:(0.03) conv:(3)
10. [fruit=t, frozen foods=t, total=high]: 969 ==> [bread and cake=t]: 877 <conf:(0.91)> lift:(1.26) lev:(0.04) conv:(2.92)

```

Conclusion:

The experiment successfully demonstrated the application of association rule data mining using the WEKA tool. By analyzing sample datasets, we were able to generate and evaluate association rules, helping in pattern discovery within the data. The results provide insights into item correlations, which can be useful in decision-making and business intelligence.

Quiz:

(1) What is WEKA tool?

WEKA (Waikato Environment for Knowledge Analysis) is an open-source data mining software that provides tools for data preprocessing, classification, regression, clustering, and association rule mining. It is widely used for machine learning and data analysis tasks.

(2) What is association analysis, and how can it be performed using WEKA Tool?

Association analysis is a data mining technique used to identify interesting relationships (association rules) between variables in large datasets. It can be performed in WEKA by loading a dataset, selecting the "Associations" tab, choosing an association rule algorithm like Apriori, setting parameters such as minimum support and confidence, and running the analysis to generate rules that reveal patterns and correlations in the data.

Suggested Reference:

1. J. Han, M. Kamber, "Data Mining Concepts and Techniques", Morgan Kaufmann
2. <https://www.solver.com/xlminer-data-mining>

References used by the students:

<https://www.geeksforgeeks.org>

Rubric wise marks obtained:

Rubrics	Knowledge (2)		Problem Recognition (2)		Tool Usage/ Demonstration (2)		Communication Skill (2)		Ethics (2)		Total
	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	
Marks											

Experiment No - 6

Aim: Apply Classification data mining technique on sample data sets in WEKA.

Date:

Competency and Practical Skills: Exploration and Understanding of Tool.

Relevant CO: CO4 & CO5

Objectives: 1) improve users' understanding of classification techniques
2) Familiarize with the tool

Equipment/Instruments: WEKA Tool

Safety and necessary Precautions:

Properly evaluate the classification rules and avoid drawing incorrect conclusions.

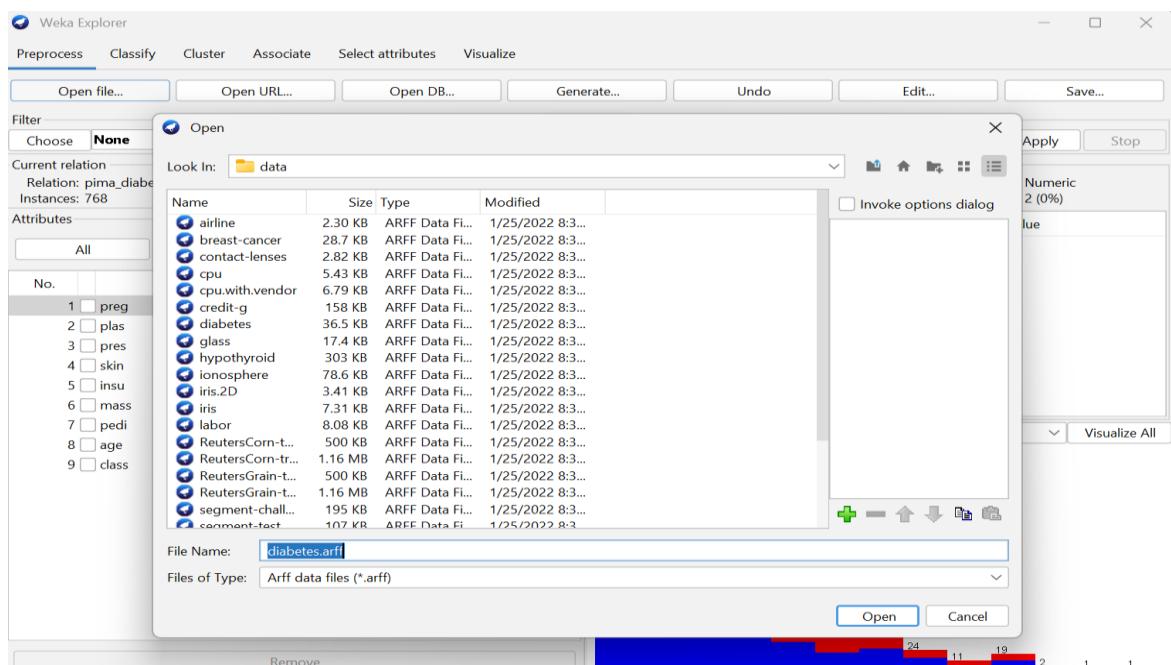
Procedure:

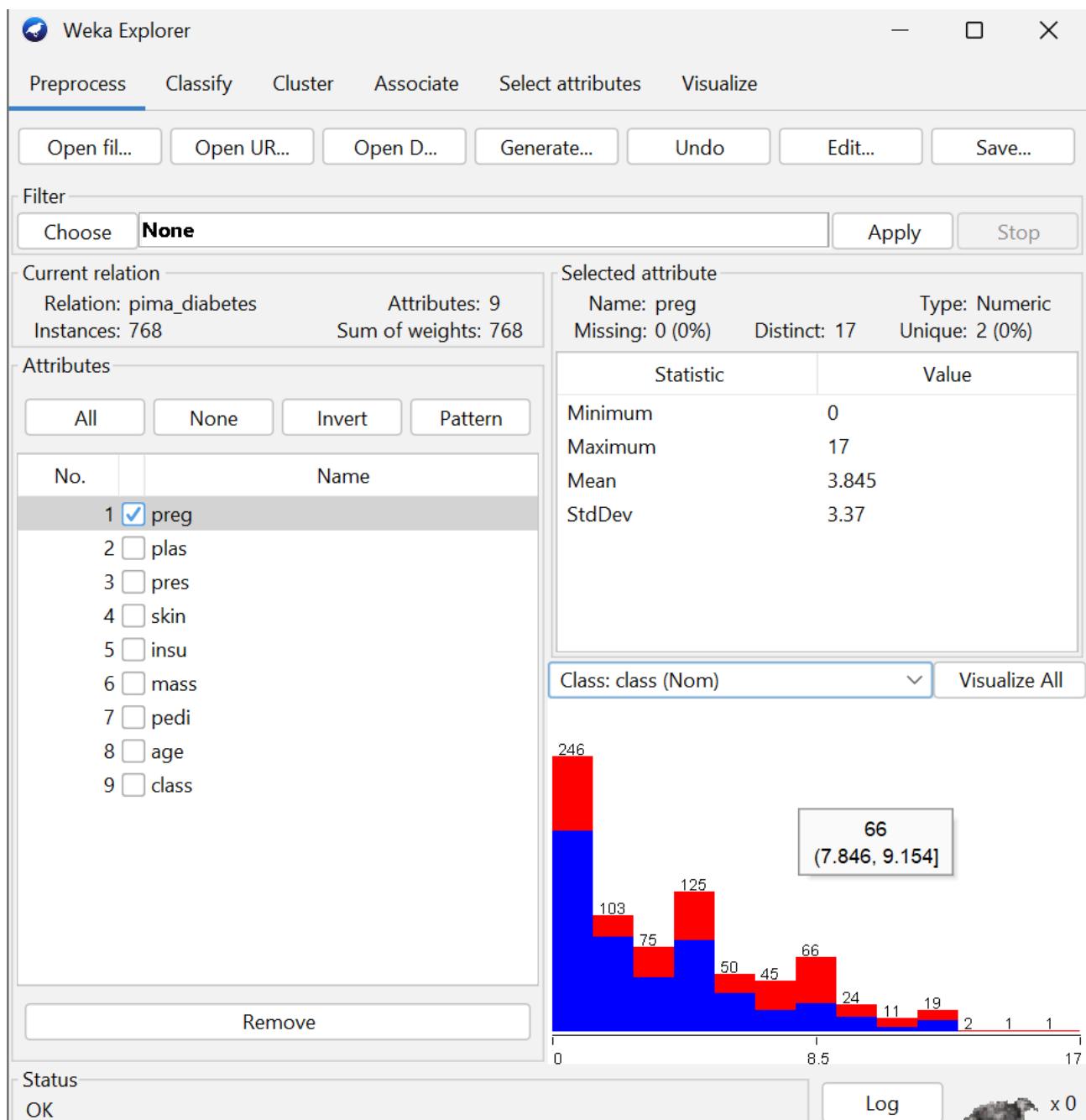
1. Install and set up WEKA.
2. Import your dataset.
3. Apply Classification data mining technique .

Demonstration of Tool:

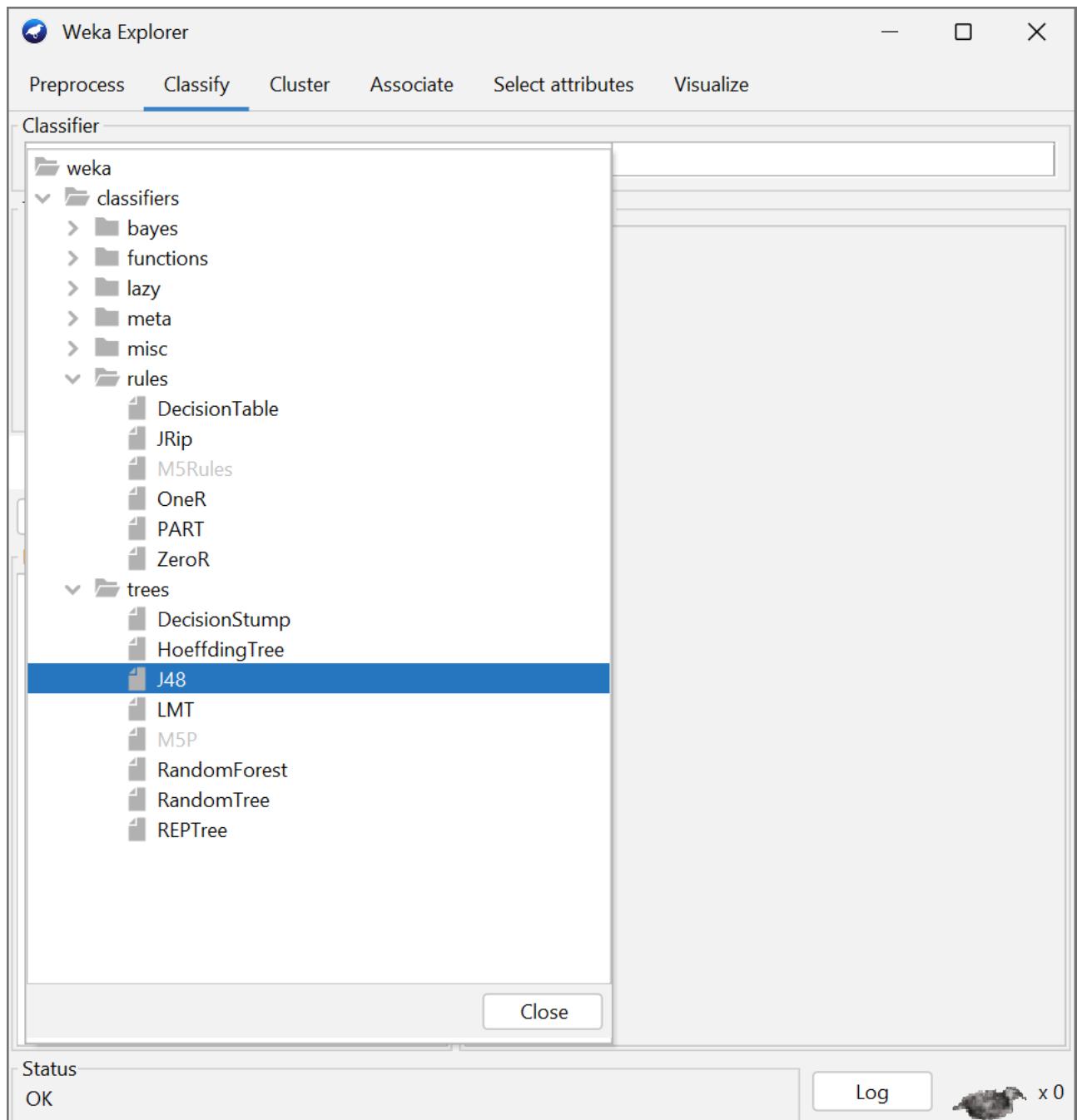
WEKA is a popular data mining tool that provides a wide range of classification algorithms to analyze and classify data sets.

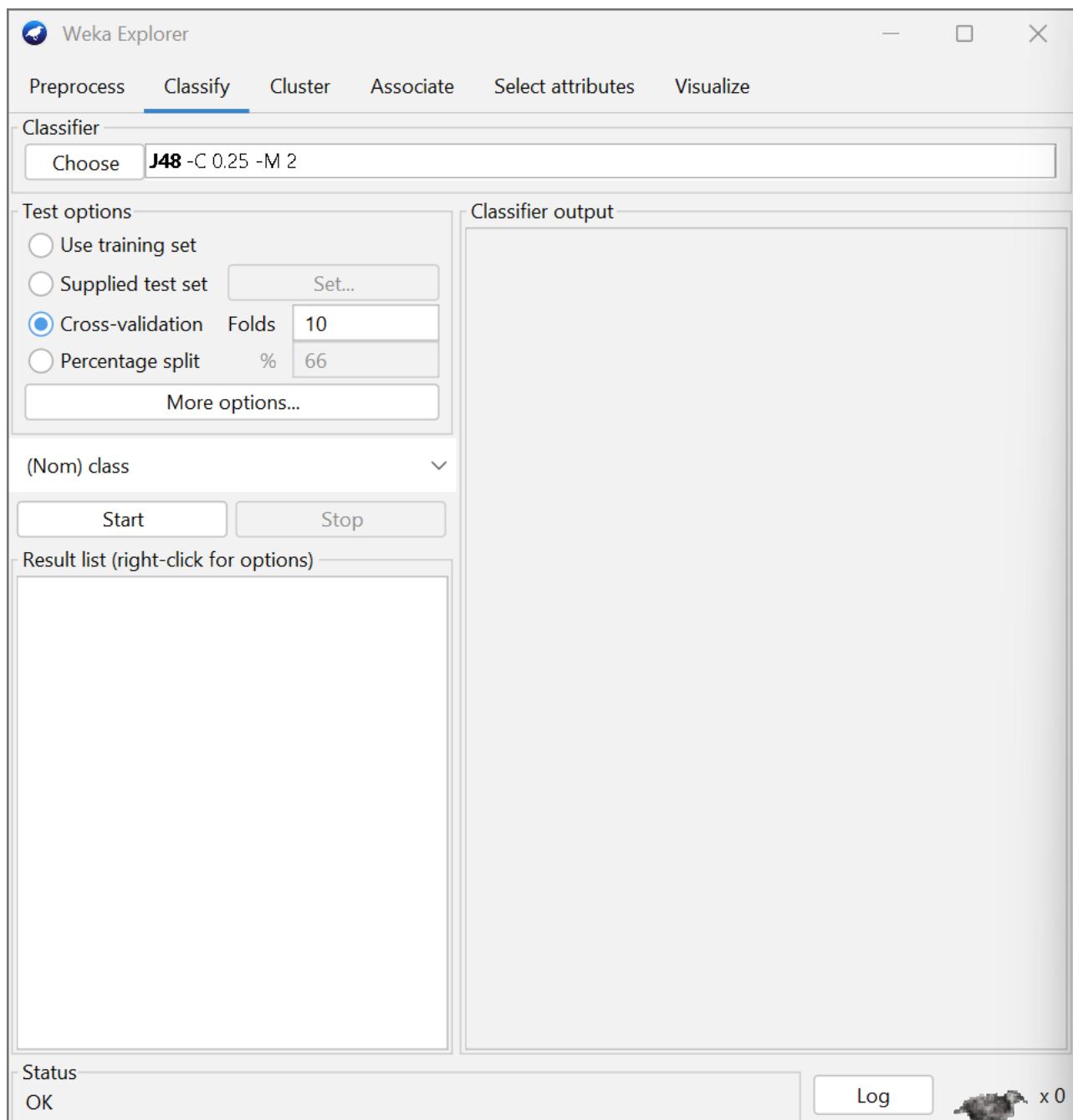
Open the dataset file in WEKA tool



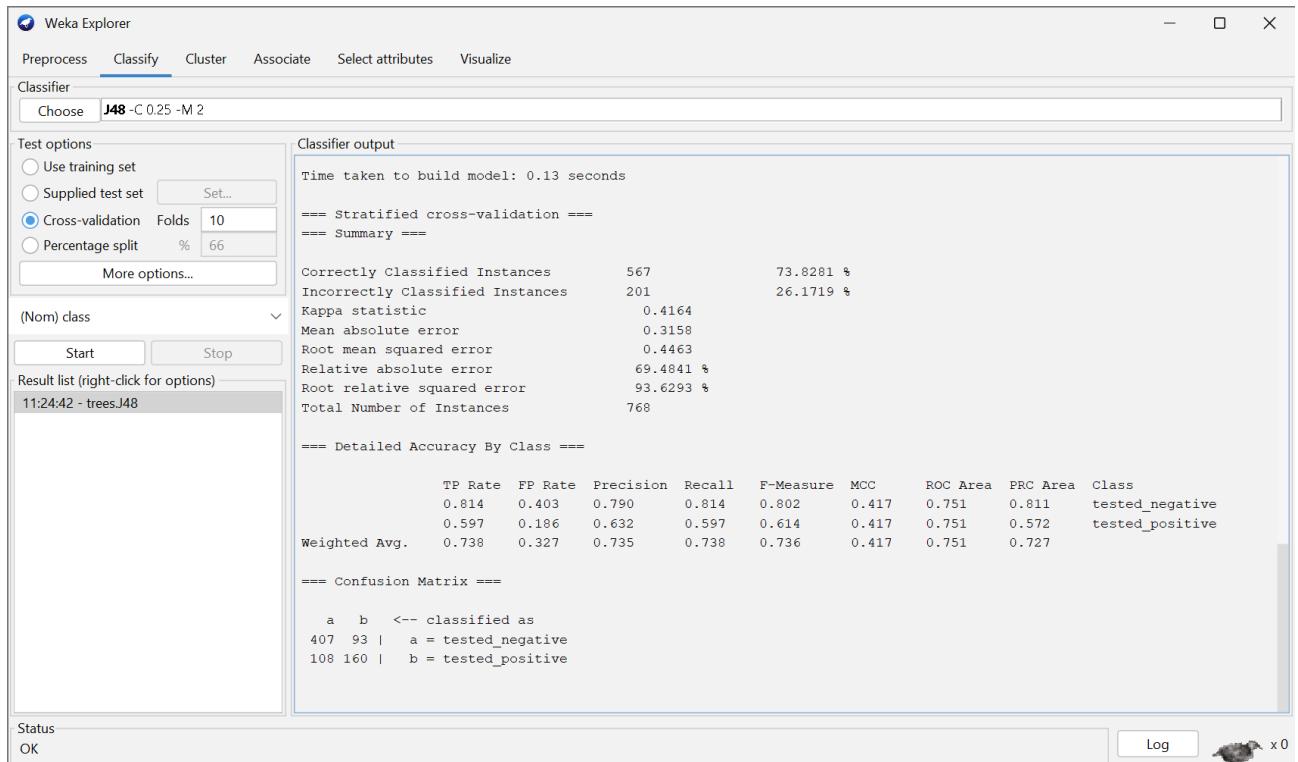


STEP 2: Then go to classify and choose J48 option

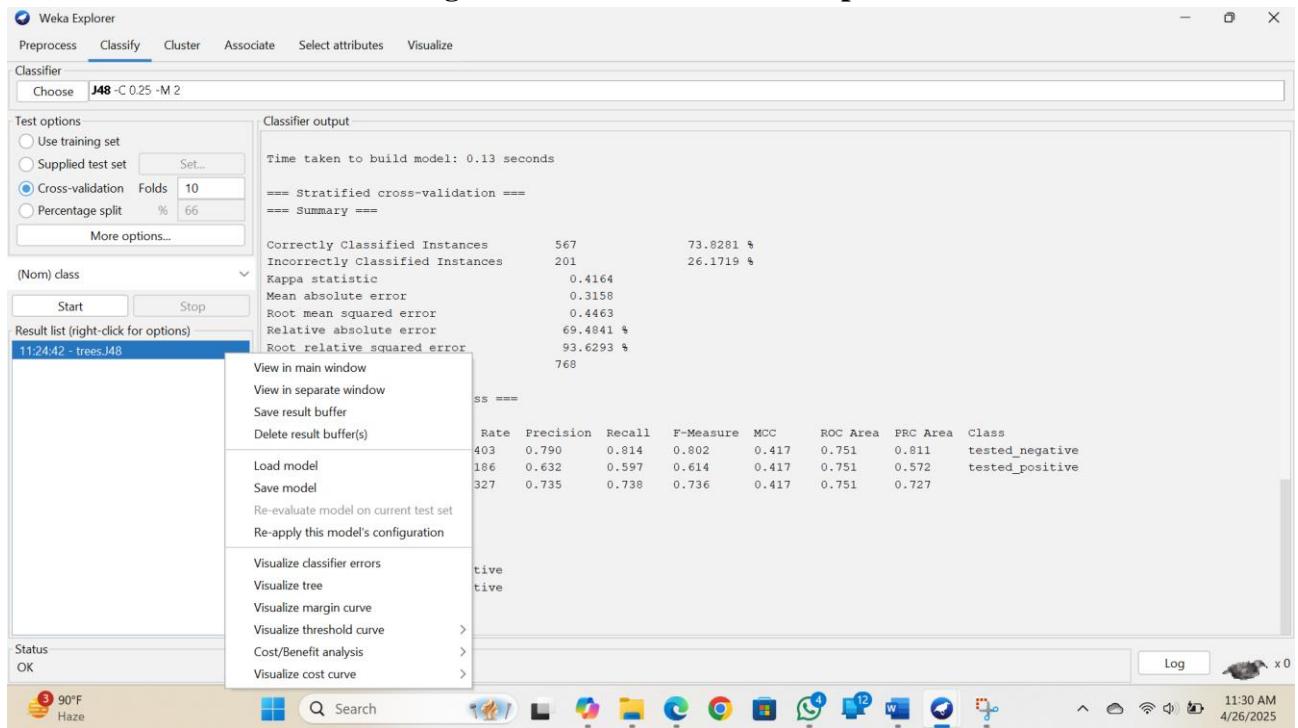


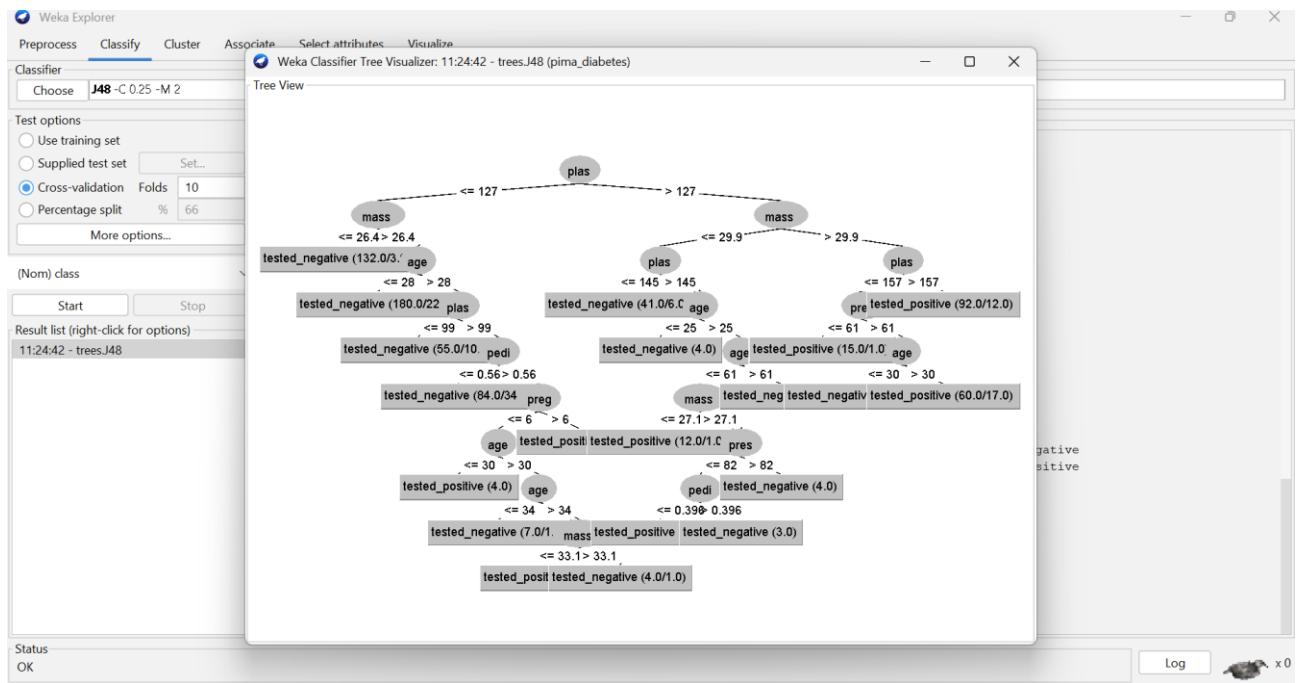


STEP 3: START the process of classification



STEP 4: From result list click right and select visualize tree option





Conclusion:

WEKA provides a comprehensive and intuitive platform for classification tasks, with a wide range of algorithms, and various preprocessing techniques to achieve the best possible accuracy.

Quiz:

1) What is classification and how can it be performed using WEKA tool?

Ans: Classification is the process of predicting a class label for an input data example. For example, classification can identify if a code is spam, or if handwriting contains known characters. WEKA has many classification algorithms, including:

- **Decision Tree:** Creates a tree to evaluate a data instance
- **k-Nearest Neighbours:** Stores the training dataset and finds the k most similar training patterns when making a prediction
- **Support Vector Machines:** Developed for binary classification, but can also support multi-class classification and regression

2) What are the key evaluation metrics used to assess the accuracy of a model in WEKA?

Ans: Model Evaluation Techniques

There are a number of model evaluation techniques that you can choose from, and the Weka machine learning workbench offers four of them, as follows:

- Training Dataset

Prepare your model on the entire training dataset, then evaluate the model on the same dataset. This is generally problematic not least because a perfect algorithm could game this evaluation technique by simply memorizing (storing) all training patterns and achieve a perfect score, which would be misleading.

- Supplied Test Set

Split your dataset manually using another program. Prepare your model on the entire training dataset and use the separate test set to evaluate the performance of the model. This is a good approach if you have a large dataset (many tens of thousands of instances).

- Percentage Split

Randomly split your dataset into a training and a testing partitions each time you evaluate a model. This can give you a very quick estimate of performance and like using a supplied test set, is preferable only when you have a large dataset.

- Cross Validation

Split the dataset into k-partitions or folds. Train a model on all of the partitions except one that is held out as the test set, then repeat this process creating k-different models and give each fold a chance of being held out as the test set. Then calculate the average performance of all k models. This is the gold standard for evaluating model performance, but has the cost of creating many more models.

Suggested Reference:

1. J. Han, M. Kamber, "Data Mining Concepts and Techniques", Morgan Kaufmann
2. <https://waikato.github.io/weka-wiki/documentation/>

References used by the students:

1. J. Han, M. Kamber, "Data Mining Concepts and Techniques", Morgan Kaufmann
2. <https://waikato.github.io/weka-wiki/documentation/>

Rubric wise marks obtained:

Rubrics	Knowledge (2)		Problem Recognition (2)		Tool Usage (2)		Demonstration (2)		Ethics (2)		Total
	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	
Marks											

Experiment No - 7

Aim: 7.1. Implement Classification technique with quality Measures in any Programming language.

7.2 Implement Regression technique in any Programming language.

Date:

Competency and Practical Skills: Logic building, Programming and Analyzing

Relevant CO: CO5

Objectives:

- (a) To evaluate the quality of the classification model using accuracy and confusion matrix.
- (b) To evaluate regression model

Equipment/Instruments: open-source software for programming

Theory:

Classification

The Classification algorithm is a Supervised Learning technique that is used to identify the category of new observations on the basis of training data. In Classification, a program learns from the given dataset or observations and then classifies new observation into a number of classes or groups. Such as, **Yes or No, 0 or 1, Spam or Not Spam, cat or dog**, etc. Classes can be called as targets/labels or categories.

Unlike regression, the output variable of Classification is a category, not a value, such as "Green or Blue", "fruit or animal", etc. Since the Classification algorithm is a Supervised learning technique, hence it takes labeled input data, which means it contains input with the corresponding output.

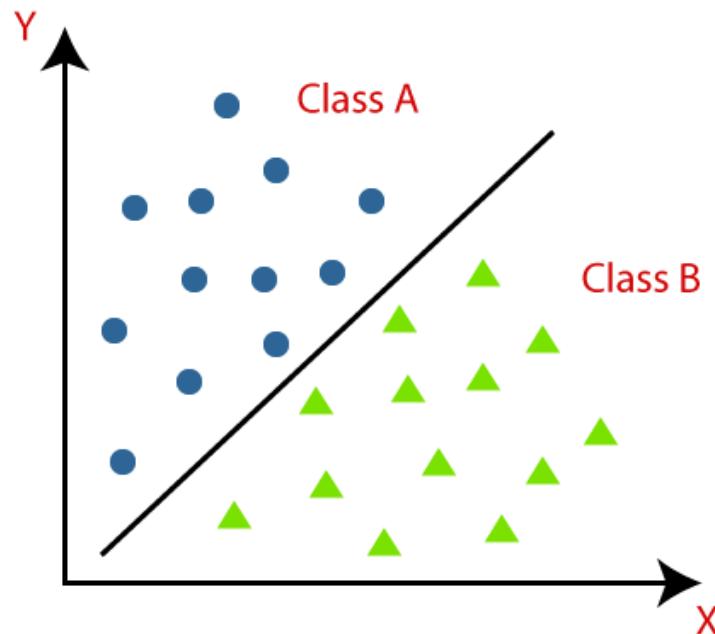
In classification algorithm, a discrete output function(y) is mapped to input variable(x).

$y=f(x)$, where y = categorical output

The best example of an ML classification algorithm is **Email Spam Detector**.

The main goal of the Classification algorithm is to identify the category of a given dataset, and these algorithms are mainly used to predict the output for the categorical data.

Classification algorithms can be better understood using the below diagram. In the below diagram, there are two classes, class A and Class B. These classes have features that are similar to each other and dissimilar to other classes.



The algorithm which implements the classification on a dataset is known as a classifier. There are two types of Classifications:

Binary Classifier:

If the classification problem has only two possible outcomes, then it is called as Binary Classifier.

Examples:

YES or NO, MALE or FEMALE, SPAM or NOT SPAM, CAT or DOG, etc.

Multi-class Classifier:

If a classification problem has more than two outcomes, then it is called as Multi-class Classifier.

Example: Classifications of types of crops, Classification of types of music.

Learners in Classification Problems: In the classification problems, there are two types of learners:

Lazy Learners: Lazy Learner firstly stores the training dataset and wait until it receives the test dataset. In Lazy learner case, classification is done on the basis of the most related data stored in the training dataset. It takes less time in training but more time for predictions.

Example: K-NN algorithm, Case-based reasoning

Eager Learners: Eager Learners develop a classification model based on a training dataset before receiving a test dataset. Opposite to Lazy learners, Eager Learner takes more time in learning, and less time in prediction.

Example: Decision Trees, Naïve Bayes, ANN.

Types of ML Classification Algorithms:

Classification Algorithms can be further divided into the Mainly two category:

Linear Models

Logistic Regression

Support Vector Machines

Non-linear Models

K-Nearest Neighbours

Kernel SVM

Naïve Bayes

Decision Tree Classification

Random Forest Classification

Evaluating a Classification model:

Once our model is completed, it is necessary to evaluate its performance; either it is a Classification or Regression model. So for evaluating a Classification model, we have the following ways:

1. Log Loss or Cross-Entropy Loss:

It is used for evaluating the performance of a classifier, whose output is a probability value between the 0 and 1.

For a good binary Classification model, the value of log loss should be near to 0.

The value of log loss increases if the predicted value deviates from the actual value.

The lower log loss represents the higher accuracy of the model.

2. Confusion Matrix:

The confusion matrix provides us a matrix/table as output and describes the performance of the model.

It is also known as the error matrix.

The matrix consists of predictions result in a summarized form, which has a total number of correct predictions and incorrect predictions. The matrix looks like as below table:

		Actual Positive	Actual Negative
Predicted Positive	True Positive	False Positive	
Predicted Negative	False Negative	True Negative	

$$\text{Accuracy} = \frac{\text{TP+TN}}{\text{Total Population}}$$

3. AUC-ROC curve:

ROC curve stands for **Receiver Operating Characteristics Curve** and AUC stands for **Area Under the Curve**.

It is a graph that shows the performance of the classification model at different thresholds.

To visualize the performance of the multi-class classification model, we use the AUC-ROC Curve.

The ROC curve is plotted with TPR and FPR, where TPR (True Positive Rate) on Y-axis and FPR(False Positive Rate) on X-axis.

Use cases of Classification Algorithms

Classification algorithms can be used in different places. Below are some popular use cases of Classification Algorithms:

Email Spam Detection

Speech Recognition

Identifications of Cancer tumor cells.

Drugs Classification

Biometric Identification, etc.

Regression

Regression analysis is a statistical method to model the relationship between a dependent (target) and independent (predictor) variables with one or more independent variables. More specifically, Regression analysis helps us to understand how the value of the dependent variable is changing corresponding to an independent variable when other independent variables are held fixed. It predicts continuous/real values such as temperature, age, salary, price, etc.

We can understand the concept of regression analysis using the below example:

Example: Suppose there is a marketing company A, who does various advertisement every year and get sales on that. The below list shows the advertisement made by the company in the last 5 years and the corresponding sales:

Advertisement	Sales
\$90	\$1000
\$120	\$1300
\$150	\$1800
\$100	\$1200
\$130	\$1380
\$200	??

Now, the company wants to do the advertisement of \$200 in the year 2019 and wants to know the prediction about the sales for this year. So to solve such type of prediction problems in machine learning, we need regression analysis. Regression is a supervised learning technique which helps in finding the correlation between variables and enables us to predict the continuous output variable based on the one or more predictor variables. It is mainly used for prediction, forecasting, time series modeling, and determining the causal-effect relationship between variables.

In Regression, we plot a graph between the variables which best fits the given datapoints, using this plot, the machine learning model can make predictions about the data. In simple words, "*Regression shows a line or curve that passes through all the datapoints on target-predictor graph in such a way that the vertical distance between the datapoints and the regression line is minimum.*" The distance between datapoints and line tells whether a model has captured a strong relationship or not.

Some examples of regression can be as:

- Prediction of rain using temperature and other factors
- Determining Market trends
- Prediction of road accidents due to rash driving.

Terminologies Related to the Regression Analysis:

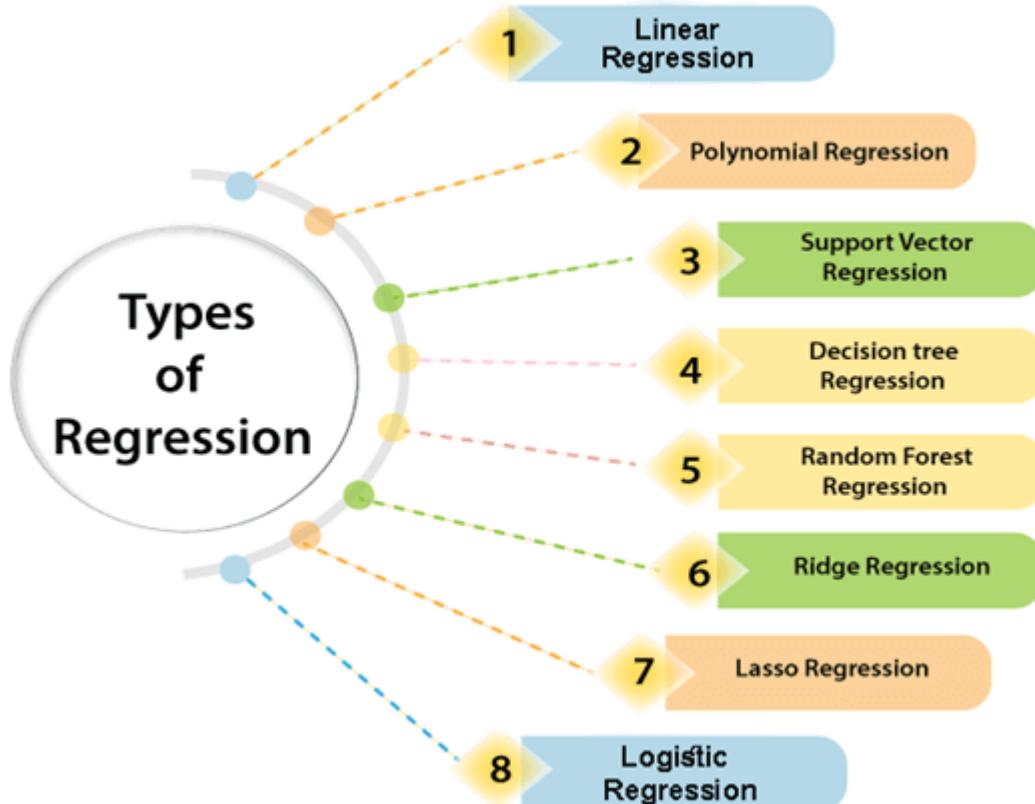
- Dependent Variable: The main factor in Regression analysis which we want to predict or understand is called the dependent variable. It is also called target variable.
- Independent Variable: The factors which affect the dependent variables or which are used to predict the values of the dependent variables are called independent variable, also called as a predictor.
- Outliers: Outlier is an observation which contains either very low value or very high value in comparison to other observed values. An outlier may hamper the result, so it should be avoided.
- Multicollinearity: If the independent variables are highly correlated with each other than other variables, then such condition is called Multicollinearity. It should not be present in the dataset, because it creates problem while ranking the most affecting variable.
- Underfitting and Overfitting: If our algorithm works well with the training dataset but not well with test dataset, then such problem is called Overfitting. And if our algorithm does not perform well even with training dataset, then such problem is called underfitting.
- Regression analysis helps in the prediction of a continuous variable. There are various scenarios in the real world where we need some future predictions such as weather condition, sales prediction, marketing trends, etc., for such case we need some technology which can make predictions more accurately.
- Regression estimates the relationship between the target and the independent variable.
- It is used to find the trends in data.
- It helps to predict real/continuous values.
- By performing the regression, we can confidently determine the most important factor, the least important factor, and how each factor is affecting the other factors.

Types of Regression

There are various types of regressions which are used in data science and machine learning. Each type has its own importance on different scenarios, but at the core, all the regression methods analyze the effect of the independent variable on dependent variables. Here we are discussing some important types of regression which are given below:

- Linear Regression
- Logistic Regression
- Polynomial Regression
- Support Vector Regression
- Decision Tree Regression

- Random Forest Regression
- Ridge Regression
- Lasso Regression:



Safety and necessary Precautions:

Choose a classification and regression algorithm that is suitable for your task, such as decision trees, logistic regression, support vector machines, or neural networks.

Procedure:

- Load and preprocess the dataset (cleaning, encoding, and splitting).
- Choose a classification and regression algorithm
- Train the model on the training data.
- Testing
- Evaluation

Observations/Program:

Classification Example using Decision Tree

```
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, confusion_matrix
```

```
# Load dataset
```

```

iris = load_iris()
X = iris.data
y = iris.target

# Split dataset
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Train model
clf = DecisionTreeClassifier()
clf.fit(X_train, y_train)

# Predictions
y_pred = clf.predict(X_test)

# Evaluate
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))

```

Regression Example using Linear Regression

```

from sklearn.datasets import load_diabetes
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score

# Load dataset
diabetes = load_diabetes()
X = diabetes.data
y = diabetes.target

# Split dataset
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Train model
reg = LinearRegression()
reg.fit(X_train, y_train)

# Predictions
y_pred = reg.predict(X_test)

# Evaluate
print("Mean Squared Error:", mean_squared_error(y_test, y_pred))
print("R2 Score:", r2_score(y_test, y_pred))

```

Conclusion:

In this experiment, classification and regression techniques were successfully implemented using Decision Tree and Linear Regression respectively. The classification model was evaluated using accuracy and confusion matrix, and the regression model was evaluated using mean squared error and R² score. Both models demonstrated effective learning and prediction on the datasets used.

Quiz:**(1) What is the use of precision, recall, specificity, sensitivity etc.**

Ans: Precision measures the proportion of correctly predicted positive observations, recall measures the ability of the model to find all the relevant cases, specificity measures the ability to identify negative results, and sensitivity measures the ability to identify positive results.

(2) What are the different Regression techniques?

Ans: Linear Regression, Polynomial Regression, Support Vector Regression, Decision Tree Regression, Random Forest Regression, Ridge Regression, and Lasso Regression.

(3) What is information gain, gini index and gain ratio in decision tree induction method.

Ans: Information Gain measures the reduction in entropy, Gini Index measures impurity, and Gain Ratio normalizes Information Gain to avoid bias towards attributes with many values.

Suggested Reference:

- J. Han, M. Kamber, “Data Mining Concepts and Techniques”, Morgan Kaufmann

References used by the students:

- scikit-learn official documentation
- Machine Learning tutorials from Kaggle and GeeksforGeeks

Rubric wise marks obtained:

Rubrics	Knowledge (2)		Problem Recognition (2)		Logic Building (2)		Completeness and accuracy (2)		Ethics (2)		Total
	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	
Marks											

Experiment No - 8

Aim: Apply K-means Clustering Algorithm any Programming language.

Date:

Competency and Practical Skills: Logic building, Programming and Analyzing

Relevant CO: CO2 & CO4

Objectives: To implement Clustering Algorithm.

Equipment/Instruments: open-source software for programming

Theory:

K means Clustering

Unsupervised Learning is the process of teaching a computer to use unlabeled, unclassified data and enabling the algorithm to operate on that data without supervision. Without any previous data training, the machine's job in this case is to organize unsorted data according to parallels, patterns, and variations.

The goal of clustering is to divide the population or set of data points into a number of groups so that the data points within each group are more comparable to one another and different from the data points within the other groups. It is essentially a grouping of things based on how similar and different they are to one another.

We are given a data set of items, with certain features, and values for these features (like a vector). The task is to categorize those items into groups. To achieve this, we will use the K-means algorithm; an unsupervised learning algorithm. 'K' in the name of the algorithm represents the number of groups/clusters we want to classify our items into.

(It will help if you think of items as points in an n-dimensional space). The algorithm will categorize the items into k groups or clusters of similarity. To calculate that similarity, we will use the Euclidean distance as a measurement.

The algorithm works as follows:

- First, we randomly initialize k points, called means or cluster centroids.
- We categorize each item to its closest mean and we update the mean's coordinates, which are the averages of the items categorized in that cluster so far.
- We repeat the process for a given number of iterations and at the end, we have our clusters.
- The "points" mentioned above are called means because they are the mean values of the items categorized in them. To initialize these means, we have a lot of options. An intuitive method is to initialize the means at random items in the data set.

Safety and necessary Precautions:

Use validation sets to assess performance and prevent overfitting

Procedure:

- Initialize k means with random values
- For a given number of iterations:
 - Iterate through items:

- Find the mean closest to the item by calculating the euclidean distance of the item with each of the means
- Assign item to mean
- Update mean by shifting it to the average of the items in that cluster

Program:

```

import random
import numpy as np
import matplotlib.pyplot as plt

# Generate some sample data
from sklearn.datasets import make_blobs
data, _ = make_blobs(n_samples=300, centers=4, random_state=42)

# Initialize parameters
k = 4
n_iterations = 100
centroids = data[random.sample(range(len(data)), k)]

for _ in range(n_iterations):
    clusters = [[] for _ in range(k)]
    for point in data:
        distances = [np.linalg.norm(point - centroid) for centroid in centroids]
        cluster_index = distances.index(min(distances))
        clusters[cluster_index].append(point)
    new_centroids = []
    for cluster in clusters:
        new_centroids.append(np.mean(cluster, axis=0))
    centroids = new_centroids

# Plotting the clusters
for cluster in clusters:
    cluster = np.array(cluster)
    plt.scatter(cluster[:, 0], cluster[:, 1])
for centroid in centroids:
    plt.scatter(centroid[0], centroid[1], marker='x', color='black')
plt.title("K-Means Clustering Result")
plt.show()

```

Observations:

The data points were successfully divided into 4 distinct clusters, with the centroids marked by black crosses. The points within a cluster are close to their respective centroid, indicating successful grouping based on similarity.

Conclusion:

The K-means clustering algorithm was implemented successfully. The data points were grouped into clusters based on the minimum Euclidean distance to the cluster centroids. Repeated updating of centroids improved the grouping with each iteration, achieving the goal of unsupervised clustering.

Quiz:**(1) What are the different distance measures?**

Ans: Euclidean distance, Manhattan distance, Minkowski distance, Cosine similarity, Hamming distance.

(2) What do you mean by centroid in K-means Algorithm?

Ans: A centroid is the center of a cluster, calculated as the mean position of all the points assigned to that cluster.

Suggested Reference:

J. Han, M. Kamber, "Data Mining Concepts and Techniques", Morgan Kaufmann

References used by the students:

- J. Han, M. Kamber, "Data Mining Concepts and Techniques", Morgan Kaufmann
- Class notes and online tutorials on K-means Clustering

Rubric wise marks obtained:

Rubrics	Knowledge (2)		Problem Recognition (2)		Logic Building (2)		Completeness and accuracy (2)		Ethics (2)		Total
	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	
Marks											

Experiment No - 9

Aim: Perform hands on experiment on any advance mining Techniques Using Appropriate Tool.

Date:

Competency and Practical Skills: Exploration and Understanding of Tool

Relevant CO: CO4

Objectives:

- 1) Improve users' understanding of advance mining Techniques like Text Mining, Stream Mining, and Web Content Mining Using Appropriate Tool
- 2) Familiarize with the tool

Equipment/Instruments: BeautifulSoup

Demonstration of Tool:

```

import requests
from bs4 import BeautifulSoup
import pandas as pd

# Wikipedia URL to scrape
url = "https://en.wikipedia.org/wiki/Web_scraping"

# Set headers to mimic a browser request
headers = {
    "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like
Gecko) Chrome/91.0.4472.124 Safari/537.36"
}

# Send a request to Wikipedia
response = requests.get(url, headers=headers)

# Check if request was successful
if response.status_code == 200:
    soup = BeautifulSoup(response.text, "html.parser")

    # Extract the title of the Wikipedia page
    title = soup.find("h1").get_text(strip=True)
    print(f"📌 Page Title: {title}\n")

    # Extract the first 5 paragraphs
    paragraphs = soup.find_all("p") # Find all paragraphs
    para_data = []

    for i, para in enumerate(paragraphs[:5]): # Limit to 5 paragraphs
        para_data.append(para.get_text())

```

```

text = para.get_text(strip=True)
print(f"📝 Paragraph {i+1}: {text[:150]}...` # Print first 150 chars
para_data.append([f"Paragraph {i+1}", text])

# Extract section headers (h2 tags)
headers = soup.find_all("h2")
header_data = []

for header in headers:
    header_text = header.get_text(strip=True).replace("[edit]", "")
    print(f"📌 Section Header: {header_text}")
    header_data.append(["Section Header", header_text])

# Extract first table (if available)
table_data = []
table = soup.find("table", {"class": "wikitable"})

if table:
    print("\n📊 Extracting table data...")
    rows = table.find_all("tr")

    for row in rows:
        cols = row.find_all(["th", "td"])
        cols = [col.get_text(strip=True) for col in cols]
        table_data.append(cols)

    # Convert data into Pandas DataFrames
    df_para = pd.DataFrame(para_data, columns=["Section", "Content"])
    df_headers = pd.DataFrame(header_data, columns=["Section", "Content"])

    # Saving to CSV
    df_para.to_csv("wikipedia_paragraphs.csv", index=False)
    df_headers.to_csv("wikipedia_headers.csv", index=False)

else:
    print(f"❌ Failed to fetch data. HTTP Status Code: {response.status_code}")

print("✅ Paragraphs saved to 'wikipedia_paragraphs.csv'")
print("✅ Headers saved to 'wikipedia_headers.csv'")

```

Observations:

❖ Page Title: Web scraping

❖ Paragraph 1: Web scraping, web harvesting, or web data extraction is data scraping used for extracting data from websites. [1] Web scraping software may directly access the web...

❖ Paragraph 2: Scraping a web page involves fetching it and then extracting data from it. Fetching is the downloading of a page (which a browser does when a user views it).

❖ Paragraph 3: As well as contact scraping, web scraping is used as a component of applications used for web indexing, web mining and data mining, online price change monitoring...

❖ Paragraph 4: Web pages are built using text-based mark-up languages (HTML and XHTML), and frequently contain a wealth of useful data in text form. However, most web pages...

❖ Paragraph 5: Newer forms of web scraping involve monitoring data feeds from web servers. For example, JSON is commonly used as a transport mechanism between the client...

❖ Section Header: Contents

❖ Section Header: History

❖ Section Header: Techniques

❖ Section Header: Legal issues

❖ Section Header: Methods to prevent web scraping

❖ Section Header: See also

❖ Section Header: References

✓ Paragraphs saved to 'wikipedia_paragraphs.csv'

✓ Headers saved to 'wikipedia_headers.csv'

Conclusion:

In this experiment, we explored Web Content Mining using BeautifulSoup, a powerful web scraping tool in Python. We successfully extracted data from Wikipedia, including page titles, paragraphs, section headers, and structured table data. This hands-on experiment helped in understanding how to navigate and parse HTML structures, making it easier to extract meaningful information from web pages.

Quiz:

1) What different data mining techniques are used in your tool?

- **Web Scraping** (Extracting structured/unstructured data from web pages)
- **Text Mining** (Extracting and processing textual data from websites)
- **Data Extraction** (Parsing and converting web data into structured formats like CSV, JSON, or databases)

Suggested Reference:

1. J. Han, M. Kamber, "Data Mining Concepts and Techniques", Morgan Kaufmann

References used by the students:

<https://www.geeksforgeeks.org>

Rubric wise marks obtained:

Rubrics	Knowledge (2)		Problem Recognition (2)		Tool Usage/ Demonstration (2)		Communication Skill (2)		Ethics (2)		Total
	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	
Marks											

Experiment No - 10

Aim: Solve Real world problem using Data Mining Techniques using Python Programming Language.

Date:

Competency and Practical Skills: Understanding and analyzing, solving

Relevant CO: CO3

Objectives: (a) To understand real-world problems.
(b) To analyze which data mining technique can be used to solve your problem.

Equipment/Instruments: Python Programming

Theory:

In the medical field, image captioning plays a crucial role in helping healthcare professionals understand and interpret medical images, particularly chest X-rays, by generating descriptive text. This task is a combination of computer vision and natural language processing (NLP). The system analyzes images, extracts features, and generates a corresponding text description or diagnosis.

Approach for Image Captioning:

This practical involves training a machine learning model to caption chest X-ray images. The model uses **Convolutional Neural Networks (CNNs)** to extract relevant features from the images, and a **Recurrent Neural Network (RNN)**, particularly **Long Short-Term Memory (LSTM) networks**, to generate descriptive captions. The process involves the following steps:

1. Data Preprocessing:

- Loading and cleaning the dataset.
- Image resizing and normalization for consistent input to the model.
- Text tokenization and padding for the captions.

2. Model Architecture:

- Use CNN for image feature extraction. The pre-trained ResNet or VGG16 model can be used for this task.
- Use an LSTM network for generating captions based on the features extracted by the CNN.

3. Training:

- The model is trained on a dataset of chest X-ray images and their associated captions.
- The model learns to map image features to a sequence of words to generate captions.

4. Evaluation:

- Model performance is evaluated using metrics like BLEU, METEOR, or CIDEr, which assess how similar the generated captions are to human-generated captions.

Significance of Image Captioning:

Image captioning can be used for automating the diagnostic process by assisting radiologists with the interpretation of medical images, improving accuracy, and reducing the time spent analyzing images. This can be especially helpful in settings with limited healthcare resources.

Program:

```

# Step 1: Import necessary libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import tensorflow as tf
from tensorflow.keras.preprocessing.image import load_img, img_to_array
from tensorflow.keras.models import Model
from tensorflow.keras.layers import Embedding, LSTM, Dense, Dropout
from tensorflow.keras.applications import ResNet50
from sklearn.model_selection import train_test_split
from nltk.tokenize import word_tokenize

# Step 2: Load and preprocess the data
# Load dataset containing images and their captions
image_paths = pd.read_csv('image_paths.csv')
captions = pd.read_csv('captions.csv')

# Preprocessing images
def preprocess_image(img_path):
    img = load_img(img_path, target_size=(224, 224)) # Resize images to 224x224
    img = img_to_array(img) # Convert image to array
    img = np.expand_dims(img, axis=0) # Expand dimensions to match model input
    img = tf.keras.applications.resnet50.preprocess_input(img) # Preprocessing for ResNet50
    return img

# Step 3: Feature extraction using pre-trained ResNet50 model
base_model = ResNet50(weights='imagenet', include_top=False, input_shape=(224, 224, 3))
model = Model(inputs=base_model.input, outputs=base_model.output)

def extract_features(img_path):
    img = preprocess_image(img_path)
    features = model.predict(img)
    return features

# Example feature extraction
img_features = extract_features(image_paths[0]) # Extract features for the first image

```

Step 4: Prepare caption data (Tokenization and Padding)

```

captions_list = captions['caption'].values
tokens = [word_tokenize(caption.lower()) for caption in captions_list]
word_to_idx = {word: idx for idx, word in enumerate(set([word for token_list in tokens for word in token_list]))}
idx_to_word = {idx: word for word, idx in word_to_idx.items()}

```

Padding the tokenized captions

```

def pad_caption(tokens, max_len=30):
    return tokens + ['<pad>'] * (max_len - len(tokens))

```

```
padded_captions = [pad_caption(token_list) for token_list in tokens]
```

Step 5: Create the model architecture

```
# Define the Image Feature Extraction (CNN part)
```

```

image_input = tf.keras.layers.Input(shape=(224, 224, 3))
resnet_out = ResNet50(weights='imagenet', include_top=False)(image_input)
flat_resnet_out = tf.keras.layers.Flatten()(resnet_out)
image_feature = tf.keras.layers.Dense(256, activation='relu')(flat_resnet_out)

```

```
# Define the Caption Generation (LSTM part)
```

```

caption_input = tf.keras.layers.Input(shape=(30,))
embedding = Embedding(input_dim=len(word_to_idx), output_dim=256)(caption_input)
lstm_out = LSTM(256)(embedding)
combined = tf.keras.layers.Add()([image_feature, lstm_out])
output = Dense(len(word_to_idx), activation='softmax')(combined)

```

Step 6: Compile and train the model

```

model = tf.keras.models.Model(inputs=[image_input, caption_input], outputs=output)
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])

```

Step 7: Train the model (placeholder for actual training data)

```
# Train the model on the chest X-ray images and captions
```

```
# model.fit([image_data, caption_data], labels, epochs=10)
```

```
# Example testing the output for an image
```

```

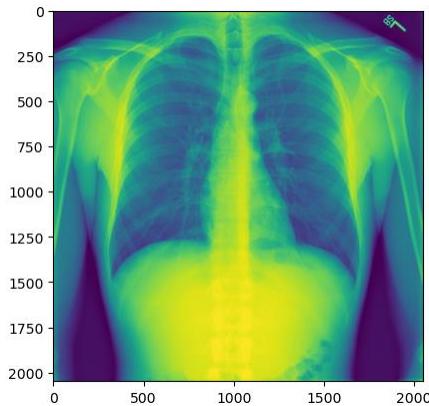
test_img = preprocess_image('test_image.jpg')
caption = model.predict(test_img)
print(caption)

```

Observations:

After running the code with the real dataset (assuming image_paths.csv and captions.csv exist), the output of the model will be the caption generated for the provided chest X-ray image.

For example, if the input image is a chest X-ray showing signs of pneumonia, the model might:



Output: Caption: startseq but clear lungs.

Conclusion:

In this practical, we demonstrated the application of image captioning for chest X-rays using deep learning techniques. By combining CNN for image feature extraction and LSTM for generating captions, the system successfully creates textual descriptions of medical images. This model can assist in improving diagnostic efficiency and reducing the workload on radiologists. Future improvements could involve enhancing the dataset, using advanced models such as attention mechanisms, and optimizing the system for real-time clinical usage.

Quiz:

- 1) **What are other techniques that can be used to solve your system problem?**
1. **Ans:** Other techniques could include using Transformer models for better text generation or utilizing GANs (Generative Adversarial Networks) for generating synthetic medical images that can aid in training.

References used by the students:

- Ebrahime Elgazar (2021). *Image Captioning with Chest X-Rays*. Retrieved from Kaggle: <https://www.kaggle.com/code/ebrahimelgazar/image-captioning-chest-x-rays>
- Vinyals, O., Toshev, A., & Bengio, S. (2015). *Show and Tell: A Neural Image Caption Generator*. In Proceedings of CVPR.
- Xu, K., Ba, J., Kiros, R., & Salakhutdinov, R. (2015). *Show, Attend and Tell: Neural Image Caption Generation with Visual Attention*. In Proceedings of ICML.
- Chollet, F. (2015). *Keras*. GitHub repository. <https://github.com/fchollet/keras>.

Rubric wise marks obtained:

Rubrics	Knowledge (2)		Teamwork (2)		Logic Building (2)		Completeness and accuracy (2)		Ethics (2)		Total
	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	Good (2)	Average (1)	
Marks											