

**INDIAN INSTITUTE OF INFORMATION TECHNOLOGY
BHOPAL**

Computer Science and Engineering Department



DATA WAREHOUSE AND MINING(DWM)(CSE-322)

PRATICAL EXAMINATION FILE

Semester-6

SUBMITTED TO:

Dr.Yatendra Sahu

SUBMITTED BY

Labhanshu Gupta

22U02007

Serial No.	Aim	Date of performance	Date of submission	Page no
1)	Data set exploration and attribute categorization	25/01/2025	29/04/2025	1-2
2)	Dataset Pre-processing and Exploration	31/01/2025	29/04/2025	3-5
3)	Bayesian Classification for a Dataset After Preprocessing using Weka and Python Libraries	3/04/2025	29/04/2025	6-11
4)	Decision Tree Classification for a Dataset used in Preprocessing using Weka and Python Libraries	3/04/2025	29/04/2025	12-15
5)	k-Nearest Neighbour Classifiers for a Dataset After Preprocessing using Weka/Python Libraries	3/04/2025	29/04/2025	16-20
6)	K-Means Clustering for a Dataset After Preprocessing using Weka/Python Libraries	8/04/2025	29/04/2025	21-24
7)	Agglomerative (bottom-up approach) Clustering for a Dataset After Preprocessing using Weka/Python Libraries	8/04/2025	29/04/2025	25-29
8)	DBSCAN (DensityBased Spatial Clustering of Applications with Noise) Clustering for a Dataset After Preprocessing using Weka/Python Libraries	8/04/2025	29/04/2025	30-36

9)	Apriori Algorithm for Association Rule Mining Using Weka/Python Libraries	25/04/2025	29/04/2025	37-42
10)	FP-Growth Algorithm for Association Rule Mining Using Weka/Python Libraries	25/04/2025	29/04/2025	43-47

Assignment-1

1) Dataset Exploration and Attribute categorization

- **Dataset Name:**Traffic Crashes Dataset
- **Source:**Kaggle
- (<https://www.kaggle.com/datasets/oktayrdeki/traffic-accidents>)
- **Domain:** The dataset contains detailed information on the traffic accidents across various regions and time periods.
- **Number of Data Objects:**8077 rows (accident records)
- **Number of features:**22 attributes,including:
 - **crash_date:** The date the accident occurred.(continuous)
 - **traffic_control_device:** The type of traffic control device involved (e.g., traffic light, sign).(nominal)
 - **weather_condition:** The weather conditions at the time of the accident.(nominal)
 - **lighting_condition:** The lighting conditions at the time of the accident.(nominal)
 - **first_crash_type:** The initial type of the crash (e.g., head-on, rear-end).(nominal)
 - **trafficway_type:** The type of roadway involved in the accident (e.g., highway, local road).(nominal)
 - **alignment:** The alignment of the road where the accident occurred (e.g., straight, curved).(nominal)
 - **roadway_surface_cond:** The condition of the roadway surface (e.g., dry, wet, icy).(nominal)
 - **road_defect:** Any defects present on the road surface.(nominal)
 - **crash_type:** The overall type of the crash.(nominal)
 - **intersection_related_i:** Whether the accident was related to an intersection.(nominal)
 - **damage:** The extent of the damage caused by the accident.(continuous)
 - **prim_contributory_cause:** The primary cause contributing to the crash.(nominal)
 - **most_severe_injury:** The most severe injury sustained in the crash.(nominal)
 - **injuries_total:** The total number of injuries reported.(continuous)
 - **injuries_fatal:** The number of fatal injuries resulting from the accident. (continuous)
 - **injuries_incapacitating:** The number of incapacitating injuries. (continuous)
 - **injuries_non_incapacitating:** The number of non-incapacitating injuries. (continuous)
 - **injuries_reported_not_evident:** The number of injuries reported but not visibly evident. (continuous)
 - **crash_hour:** The hour the accident occurred. (continuous)
 - **crash_day_of_week:** The day of the week the accident occurred.(continuous)
 - **crash_month:** The month the accident occurred.(continuous)

- **Categorial (Nominal) Attributes:**

- Traffic_control_device
- Weather_condition
- Lightining_condition
- First_crash_type
- Trafficway_type
- Alignment
- Roadway_surface_condion
- Road_defect
- Crash_type
- Intersection_related_i
- Prim_contributory_cause
- Most_severe_injury

- Real-Valued(Continuous) Attributes:

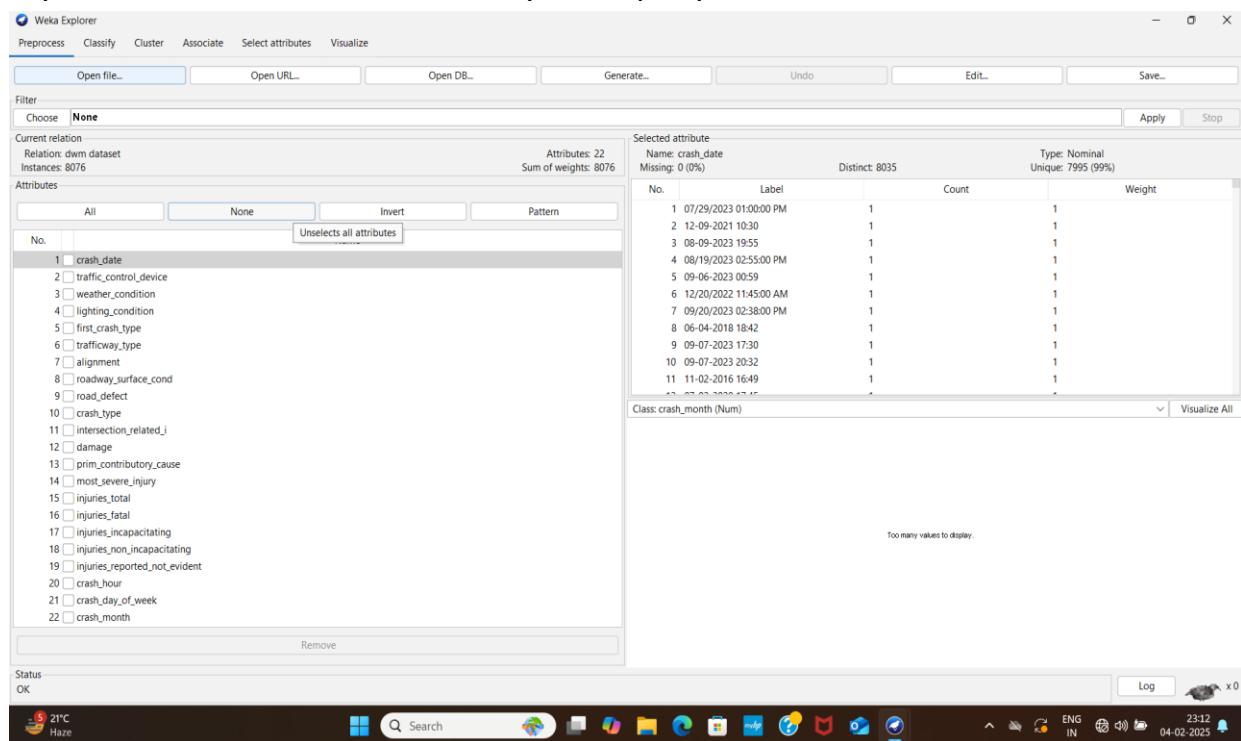
- Crash_date
- Damage
- Injuries_total
- Injuries_fatal
- Injuries_incapacitating
- injuries_non_incapacitating
- injuries_reported_not_evident
- crash_hour
- crash_day_of_week
- crash_month

Assignment-2

1)Dataset pre-processing and exploration.

1. Loading the dataset

- >Open Weka Explorer.
- >In the preprocess bar click on “open file...”.
- >select the file you want to open and press ok.
- >your data set is loaded and ready to be pre processed.



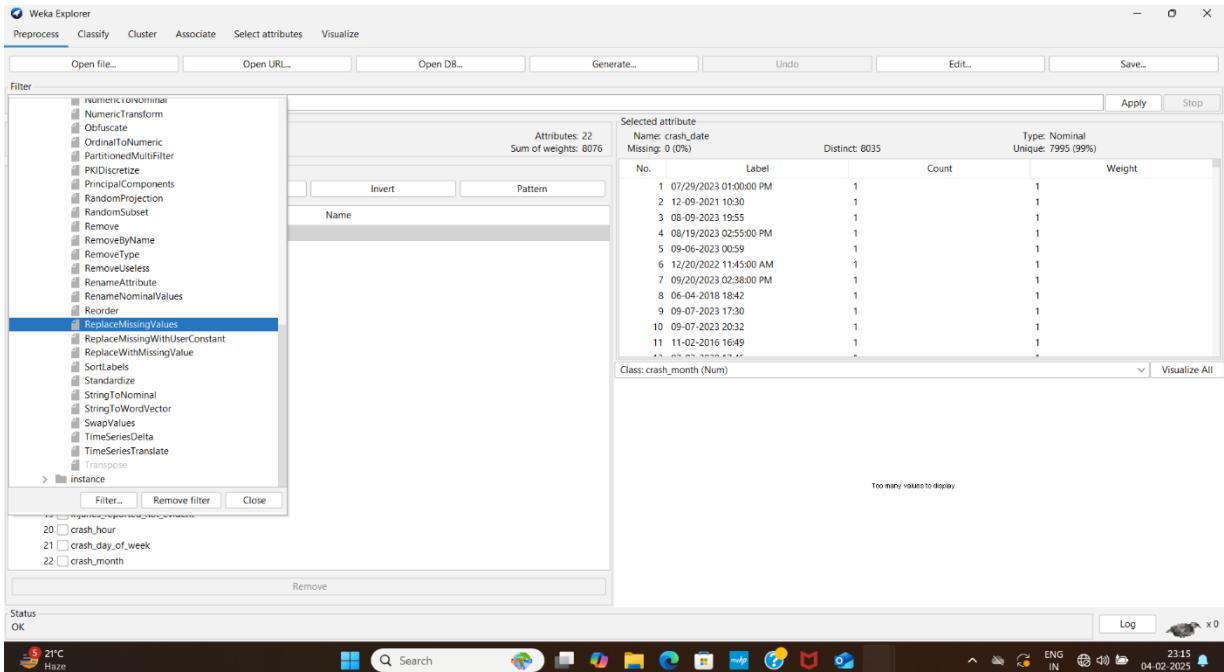
Result-> your data set file is opened in the weka explorer.

2. Handle missing values

All values are present so we don't need to apply any filter for replacing missing values

But if there were we can do it from the filters option.

- >Click on the choose option and then click on the filter option :Unsupervised
- >attribute->ReplaceMissingValues.

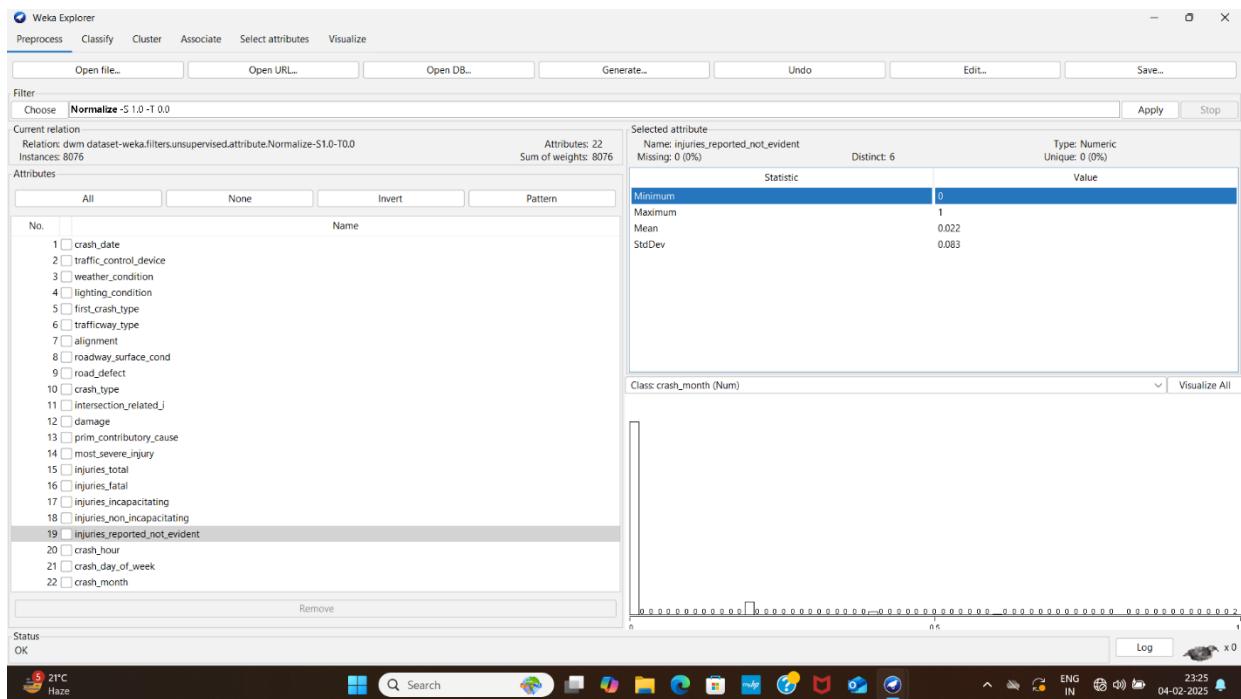


Result->This replaces the missing values with either the mean(for numeric attributes) or the mode(for nominal attributes).

3. Normalize data

->Go to the filters section again .

->Choose the filter :Unsupervised->Attribute->Normalize.



Result->All the numeric attributes are now scaled between 0 and 1.

3. Convert data types.

->Use the filter unsupervised->Attribute->NumericToNominal to convert any attribute from Numeric to normal.

->Use the filter unsupervised->Attribute->NominalToNumeric to convert any attribute from Numeric to normal.

5 Save the preprocessed dataset

->After performing all the operations Click on the save button at the top right corner of the Screen.

Result->Normalized dataset

The screenshot shows a Microsoft Excel spreadsheet titled 'sem-6.csv'. The table has 26 columns labeled A through Z. Column A contains the header 'crash_date'. The data consists of approximately 26 rows of crash information. The first few rows include:

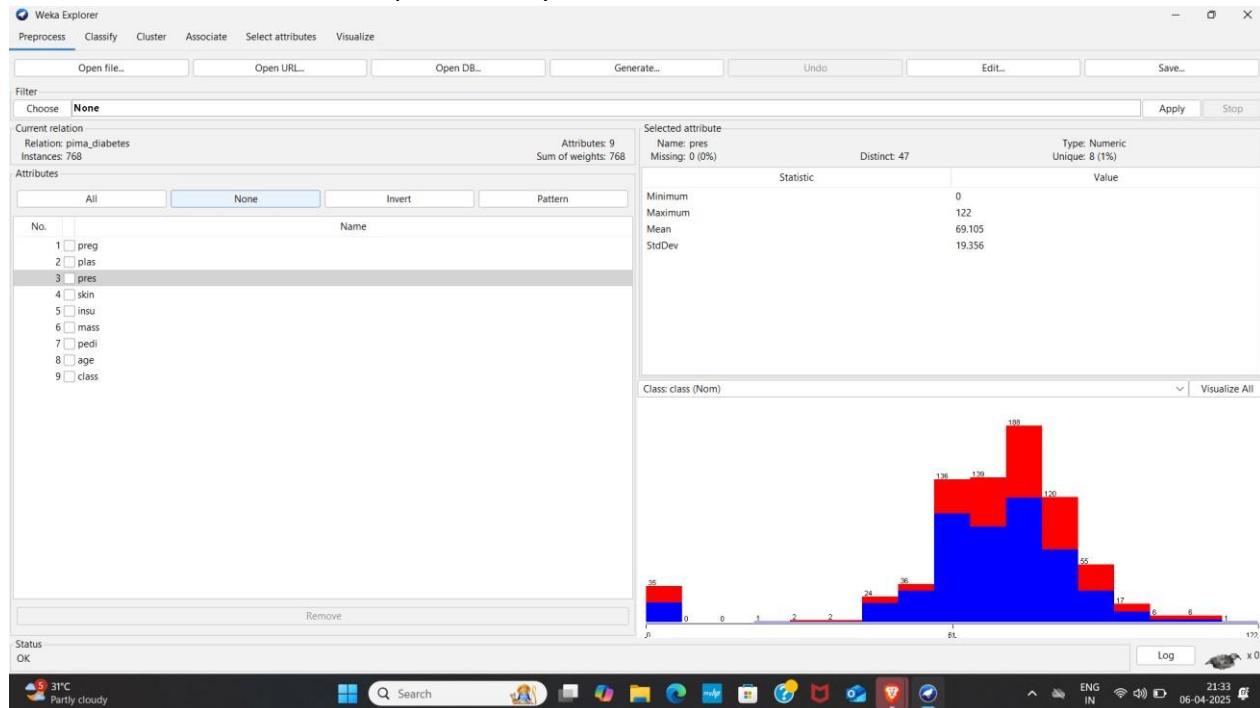
crash_date	traffic_cor	weather_c	lighting_cc	first_crash	trafficway_alignment	roadway_road	defe_crash	type_crash	intersec_damage	prim_cont	most_sev_injuries	to_injuries_fa	injuries_in	injuries_nc	injuries_re	crash_hou	crash_day	crash_month
07/29/20	TRAFFIC S CLEAR	DAYLIGHT	TURNING	'NOT DIVI'	STRAIGHT UNKNOWN	UNKNOWN	'NO INJUR Y	'\$501 - \$1 500'	'UNABLE T NO INDIC	0	0	0	0	0	13	1	7	
12-09-202	TRAFFIC S CLEAR	DAYLIGHT	'REAR ENC T-INTERSE	STRAIGHT DRY	'NO DEFEC	'NO INJUR Y	'\$501 - \$1 500'	'FOLLOWII NO INDIC	0	0	0	0	0	10	0.666667	12		
08-09-202	TRAFFIC S CLEAR	DAYLIGHT	ANGLE	'FOUR WA'	STRAIGHT DRY	'NO DEFEC	'INJURY AIY	'OVER \$1 500'	'UNABLE T NONINCA 0.833333	0	0	0.833333	0	0	19	0.5	8	
08/19/20	TRAFFIC S CLEAR	DAYLIGHT	'REAR ENC T-INTERSE	STRAIGHT UNKNOWN	UNKNOWN	'NO INJUR Y	'\$501 - \$1 500'	'DRIVING 'NO INDIC	0	0	0	0	0	14	1	8		
09-06-202	'NO CONT RAIN	'DARKNES	LIGHTED	I FIXED OBI	'NOT DIVI'	STRAIGHT WET	UNKNOWN	'INJURY AI N	'\$501 - \$1 500'	'UNABLE T NONINCA 0.333333	0	0	0.333333	0	0	0	0.5	

Assignment-3

1) Bayesian classification for a dataset after pre processing using weka and python libraries

1. Loading the dataset

Diabetes dataset is chosen to perform Bayesian classification.

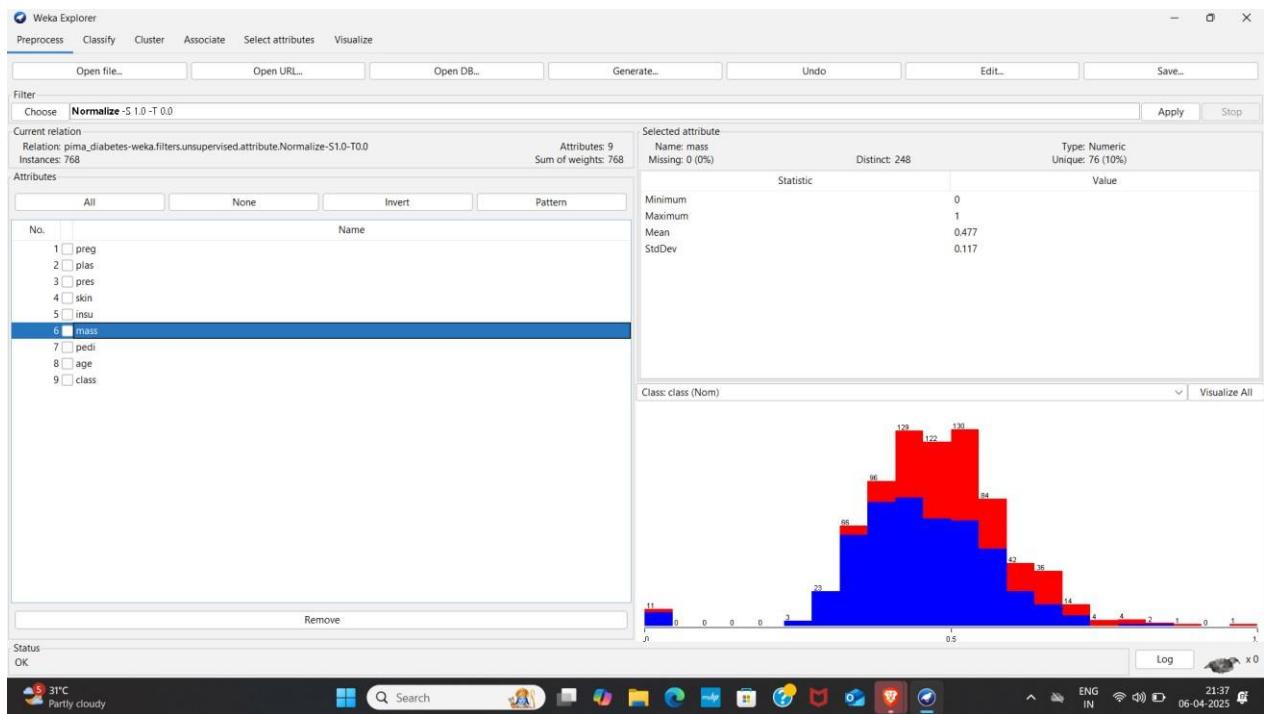


2. Handle missing values

All values are present so we don't need to apply any filter for replacing missing values

3. Normalize data

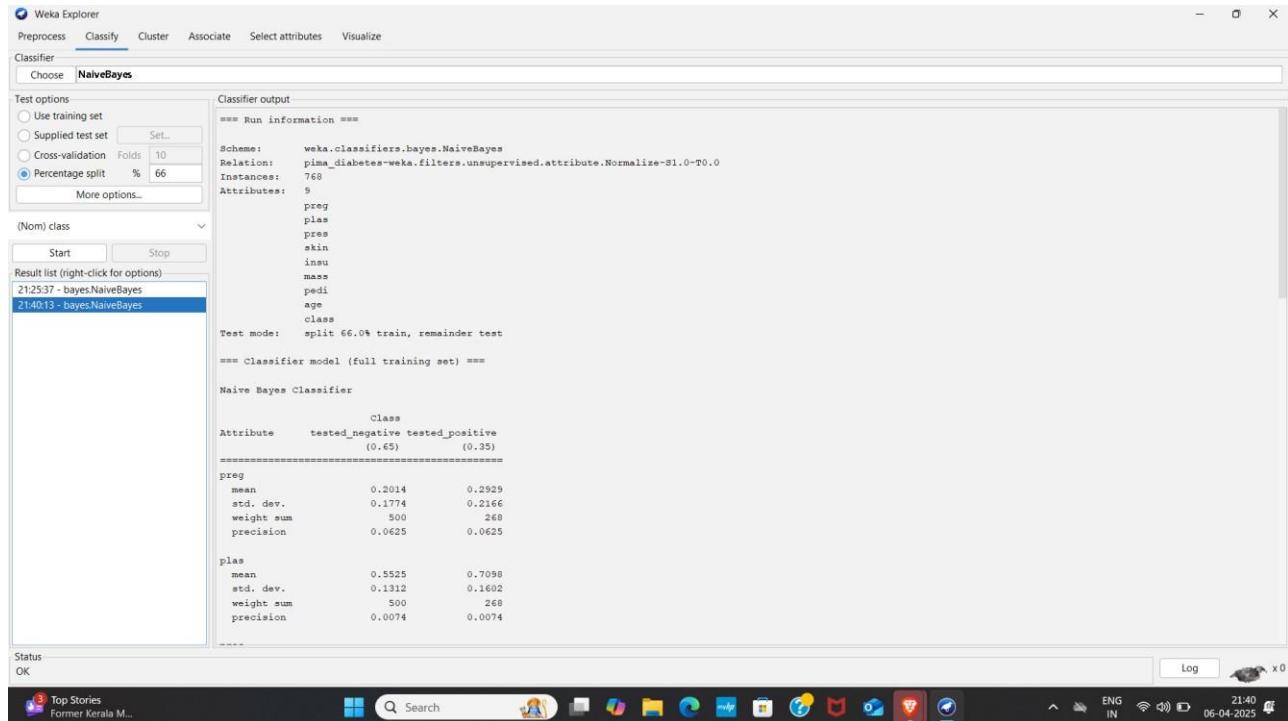
- >Go to the filters section again .
- >Choose the filter :Unsupervised->Attribute->Normalize.



Result->All the numeric attributes are now scaled between 0 and 1.

4)Applying Bayesian classification on the preprocessed data.

- 1) Click on classify in the tab above.
- 2)Choose the classifier->Bayes->NaiveBayes.
- 3)choose the values of cross-validation and percentage split
- 4)click on start.



```

Time taken to build model: 0 seconds

==== Evaluation on test split ====

Time taken to test model on test split: 0.01 seconds

==== Summary ====
Correctly Classified Instances      203      77.7778 %
Incorrectly Classified Instances   58       22.2222 %
Kappa statistic                   0.4776
Mean absolute error               0.2654
Root mean squared error          0.382
Relative absolute error          58.8383 %
Root relative squared error     81.6152 %
Total Number of Instances        261

==== Detailed Accuracy By Class ====

```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.854	0.386	0.826	0.854	0.840	0.478	0.853	0.919	0.919	tested_negative
0.614	0.146	0.662	0.614	0.638	0.478	0.853	0.760	0.760	tested_positive
Weighted Avg.	0.778	0.309	0.774	0.778	0.775	0.478	0.853	0.868	

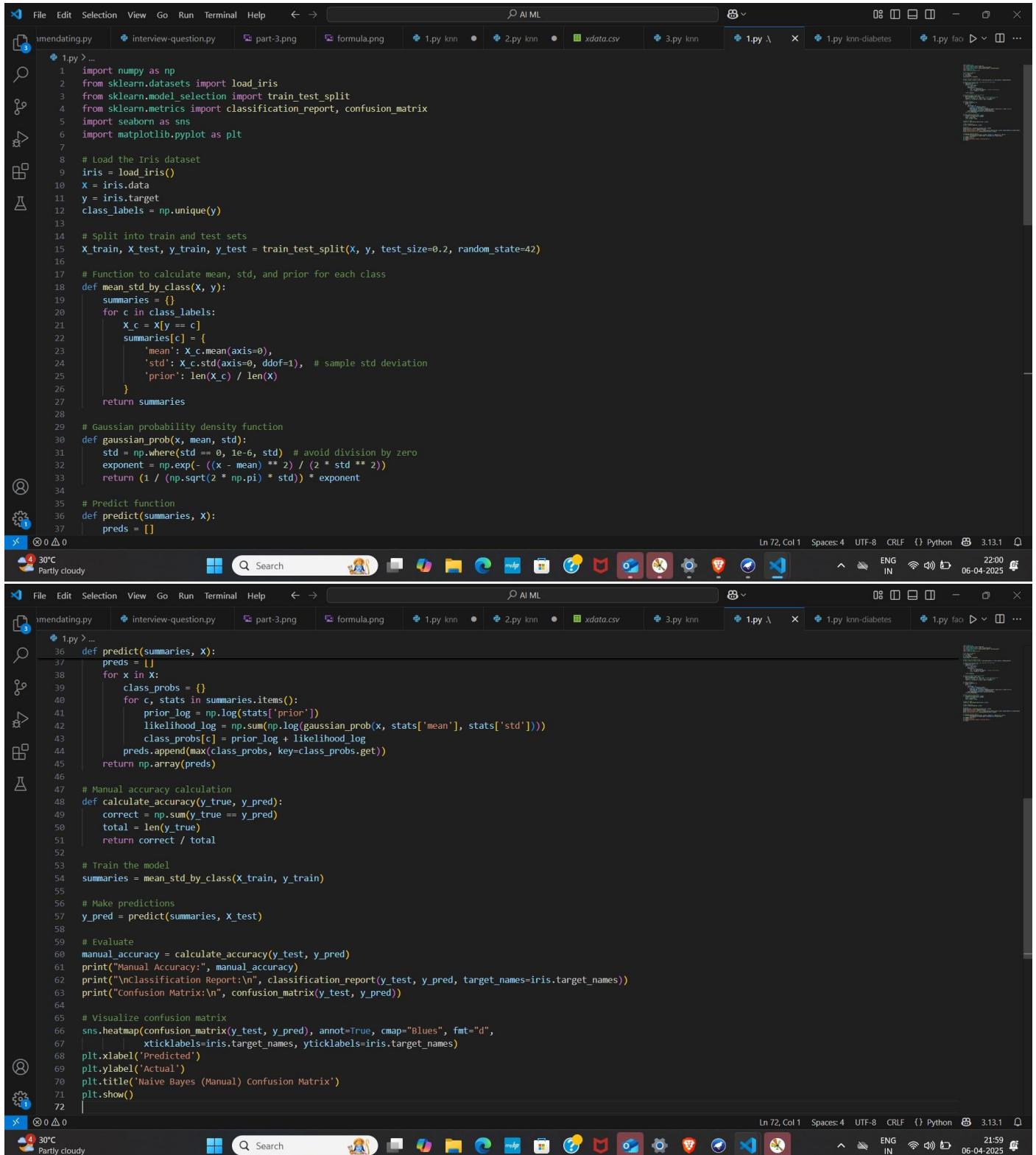
```

==== Confusion Matrix ====
a   b   <-- classified as
152 26 |   a = tested_negative
32 51 |   b = tested_positive

```

Result->Naïve's baise classification is performed and accuracy of this model is 77.778 %.

2)Python implementation for Bayesian classification



```
1 import numpy as np
2 from sklearn.datasets import load_iris
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import classification_report, confusion_matrix
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7
8 # Load the Iris dataset
9 iris = load_iris()
10 X = iris.data
11 y = iris.target
12 class_labels = np.unique(y)
13
14 # split into train and test sets
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
16
17 # Function to calculate mean, std, and prior for each class
18 def mean_std_by_class(X, y):
19     summaries = {}
20     for c in class_labels:
21         X_c = X[y == c]
22         summaries[c] = {
23             'mean': X_c.mean(axis=0),
24             'std': X_c.std(axis=0, ddof=1), # sample std deviation
25             'prior': len(X_c) / len(X)
26         }
27     return summaries
28
29 # Gaussian probability density function
30 def gaussian_prob(x, mean, std):
31     std = np.where(std == 0, 1e-6, std) # avoid division by zero
32     exponent = np.exp(-(np.sum((x - mean) ** 2) / (2 * std ** 2)))
33     return (1 / (np.sqrt(2 * np.pi) * std)) * exponent
34
35 # Predict function
36 def predict(summaries, X):
37     preds = []
38
39     for x in X:
40         class_probs = {}
41         for c, stats in summaries.items():
42             prior_log = np.log(stats['prior'])
43             likelihood_log = np.sum(np.log(gaussian_prob(x, stats['mean'], stats['std'])))
44             class_probs[c] = prior_log + likelihood_log
45         preds.append(max(class_probs, key=class_probs.get))
46
47     return np.array(preds)
48
49 # Manual accuracy calculation
50 def calculate_accuracy(y_true, y_pred):
51     correct = np.sum(y_true == y_pred)
52     total = len(y_true)
53     return correct / total
54
55 # Train the model
56 summaries = mean_std_by_class(X_train, y_train)
57
58 # Make predictions
59 y_pred = predict(summaries, X_test)
60
61 # Evaluate
62 manual_accuracy = calculate_accuracy(y_test, y_pred)
63 print("Manual Accuracy:", manual_accuracy)
64 print("\nClassification Report:\n", classification_report(y_test, y_pred, target_names=iris.target_names))
65 print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))
66
67 # Visualize confusion matrix
68 sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, cmap="Blues", fmt="d",
69             xticklabels=iris.target_names, yticklabels=iris.target_names)
70 plt.xlabel('Predicted')
71 plt.ylabel('Actual')
72 plt.title('Naive Bayes (Manual) Confusion Matrix')
73 plt.show()
```

Result:-As the dataset was clean and well seperated we got accuracy 1.0.

```

Manual Accuracy: 1.0

Classification Report:
precision    recall   f1-score   support
setosa       1.00      1.00      1.00      50
versicolor   1.00      1.00      1.00       9
virginica    1.00      1.00      1.00      11

accuracy          1.00      1.00      1.00      50
macro avg       1.00      1.00      1.00      50
weighted avg    1.00      1.00      1.00      50

Confusion Matrix:
[[10  0  0]
 [ 0  9  0]
 [ 0  0 11]]

accuracy          1.00      1.00      1.00      50
macro avg       1.00      1.00      1.00      50
weighted avg    1.00      1.00      1.00      50

Confusion Matrix:
[[10  0  0]
 [ 0  9  0]
 [ 0  0 11]]
[[10  0  0]
 [ 0  9  0]
 [ 0  0 11]]
```

3)Advantages and disadvantages of Bayesian classification is:-

Advantages:-

1) Simple and Fast:

- Easy to implement and computationally efficient.
- Works well with large datasets and high-dimensional data.

2) Performs Well with Less Data:

- Requires less training data compared to other classifiers.
- Handles both binary and multiclass classification.

3) Effective for Certain Problems:

- Surprisingly effective for text classification (e.g., spam detection, sentiment analysis).
- Performs well when the assumption of feature independence holds.

4) Handles Missing Data:

- Can handle missing values relatively well by simply ignoring the feature during probability calculation.

5) Scalable:

- Scales linearly with the number of features and data points.

Disadvantages:-

1) Strong Independence Assumption:

- Assumes features are independent, which is rarely true in real-world data.
- Violations of this assumption can lead to poor performance.

2) Zero-Frequency Problem:

- If a categorical variable has a category in the test data not seen in the training data, it assigns zero probability.
- This can be mitigated with techniques like Laplace smoothing.

3) Not Ideal for Continuous Features (without modification):

- Basic Naive Bayes works better with categorical input.
- For continuous features, it assumes a Gaussian distribution which might not always hold true.

4) Less Flexible Compared to Other Models:

- Doesn't learn interactions between features.
- Other models like decision trees or random forests can model more complex relationships.

5) Can Be Outperformed by More Complex Models:

- In practice, more sophisticated algorithms (e.g., SVM, XGBoost, neural networks) often yield better results, especially for complex datasets.

Assignment-4

1) Decision tree classification for a dataset used in preprocessing using weka and python libraries.

python implementation for decision tree classification

1) We have used titanic data set to perform decision tree classification on

Features-> PassengerId,Survived,Pclass,Name,Sex,Age,SibSp,Parch,Ticket,Fare,Cabin,Embarked.

Label->Survived .

Features we have selected are:- Pclass, Sex, Age, Fare, Embarked.

```
File Edit Selection View Go Run Terminal Help <- > AI ML 08 □ □ □ ... decision-tree.py x titanic.csv 1.py kmn 2.py kmn 3.py kmn 1.py \ 1.py kmn-diabet > ... decision-tree.py > ... 1 import csv 2 import math 3 import random 4 from collections import Counter 5 6 # Load CSV 7 def load_csv(filename): 8     with open(filename, 'r') as f: 9         reader = csv.DictReader(f) 10        data = [] 11        for row in reader: 12            if row['Age'] and row['Fare']: # skip rows with missing values 13                data.append({ 14                    'Pclass': int(row['Pclass']), 15                    'Sex': row['Sex'], 16                    'Age': float(row['Age']), 17                    'Fare': float(row['Fare']), 18                    'Embarked': row['Embarked'], 19                    'Survived': row['Survived'] 20                }) 21        return data 22 23 # Encode categorical variables manually 24 def encode(data): 25    sex_map = {'male': 0, 'female': 1} 26    embark_map = {'S': 0, 'C': 1, 'Q': 2} 27    for row in data: 28        row['Sex'] = sex_map.get(row['Sex'], -1) 29        row['Embarked'] = embark_map.get(row['Embarked'], -1) 30    return data 31 32 # Convert to list format 33 def to_row_list(data): 34    return [[row['Pclass'], row['Sex'], row['Age'], row['Fare'], row['Embarked'], row['Survived']] for row in data] 35 36 # Entropy 37 def entropy(data): 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73
```

This screenshot shows the first part of a Python script for a decision tree. It includes imports for csv, math, and random, along with a Counter from collections. The `load_csv` function reads a CSV file, skipping rows with missing age or fare values. The `encode` function maps categorical variables ('Sex' and 'Embarked') to numerical values (-1 for missing). The `to_row_list` function converts the data into a list of lists for easier processing. The `entropy` function is defined but not yet implemented.

```
File Edit Selection View Go Run Terminal Help <- > AI ML 08 □ □ □ ... decision-tree.py x titanic.csv 1.py kmn 2.py kmn 3.py kmn 1.py \ 1.py kmn-diabet > ... decision-tree.py > ... 37 def entropy(data): 38    label_counts = Counter([row[-1] for row in data]) 39    total = len(data) 40    return -sum((c/total) * math.log2(c/total) for c in label_counts.values()) 41 42 # Split data 43 def split_data(data, col, value): 44    left, right = [], [] 45    for row in data: 46        if isinstance(value, float) or isinstance(value, int): 47            if (left if row[col] <= value else right).append(row) 48        else: 49            (left if row[col] == value else right).append(row) 50    return left, right 51 52 # Best split 53 def best_split(data): 54    base_entropy = entropy(data) 55    best_col, best_val = 0, None, None 56    n_features = len(data[0]) - 1 57    for col in range(n_features): 58        values = set(row[col] for row in data) 59        for val in values: 60            left, right = split_data(data, col, val) 61            if not left or not right: 62                continue 63            p = len(left) / len(data) 64            gain = base_entropy - (p * entropy(left) + (1 - p) * entropy(right)) 65            if gain > best_gain: 66                best_gain, best_col, best_val = gain, col, val 67    return best_col, best_val 68 69 # Build tree 70 def build_tree(data): 71    labels = [row[-1] for row in data] 72    if labels.count(labels[0]) == len(labels): 73        return labels[0]
```

This screenshot shows the completed decision tree implementation. It includes the `entropy` function, which calculates the entropy of a dataset based on the frequency of each class. The `split_data` function splits the data into two sets based on a specific column and value. The `best_split` function finds the best split by calculating the information gain for each feature and returning the one with the highest gain. The `build_tree` function is defined but not yet implemented.

```
File Edit Selection View Go Run Terminal Help AI ML decision-tree.py X titanic.csv 1.py knn 2.py knn 3.py knn 1.py λ ... latint.py interview-question.py part-3.png formula.png 1.py knn 2.py knn 3.py knn 1.py λ ... 70 def build_tree(data): 71     if len(data) == 1: 72         return labels[0] 73     col, val = best_split(data) 74     if col is None: 75         return Counter(labels).most_common(1)[0][0] 76     left, right = split_data(data, col, val) 77     return { 78         'col': col, 79         'val': val, 80         'left': build_tree(left), 81         'right': build_tree(right) 82     } 83 84 85 # Predict 86 def predict(tree, row): 87     if isinstance(tree, str): 88         return tree 89     val = row[tree['col']] 90     branch = tree['left'] if (val <= tree['val']) else tree['right'] 91     return predict(branch, row) 92 93 # Accuracy 94 def accuracy(tree, data): 95     correct = 0 96     for row in data: 97         if predict(tree, row) == row[-1]: 98             correct += 1 99     return correct / len(data) 100 101 # Split data 102 def train_test_split(data, test_ratio=0.2): 103     random.shuffle(data) 104     split = int(len(data) * (1 - test_ratio)) 105     return data[:split], data[split:] 106 107 # Main run 108 data = load_csv('titanic.csv') 109 110 111 112 113 114 115 116 117
```

```
File Edit Selection View Go Run Terminal Help AI ML decision-tree.py X titanic.csv 1.py knn 2.py knn 3.py knn 1.py λ ... latint.py interview-question.py part-3.png formula.png 1.py knn 2.py knn 3.py knn 1.py λ ... 108 data = load_csv('titanic.csv') 109 encoded = encode(data) 110 rows = to_row_list(encoded) 111 112 train_data, test_data = train_test_split(rows) 113 tree = build_tree(train_data) 114 acc = accuracy(tree, test_data) 115 116 print("Custom Decision Tree Accuracy on Titanic:", round(acc * 100, 2), "%") 117
```

Result:-Accuracy we got is 77.62%.

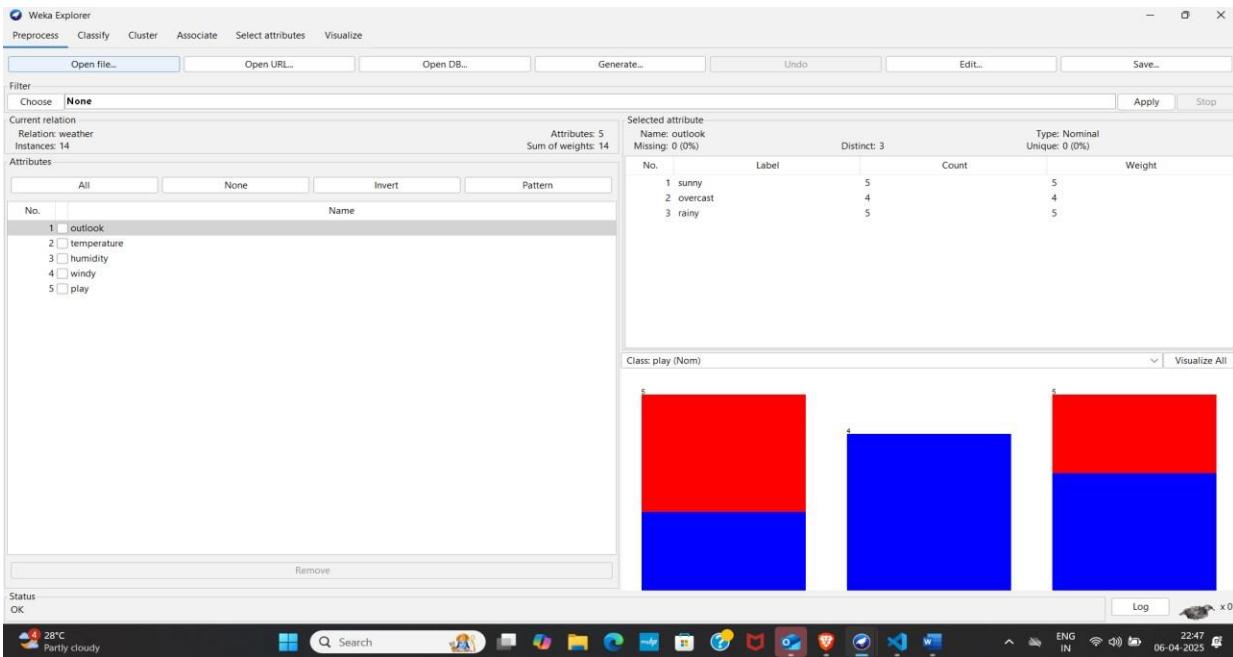
```
PS D:\AI ML> & "C:/Program Files/Python313/python.exe" "d:/AI ML/decision-tree.py"
Custom Decision Tree Accuracy on Titanic: 77.62 %
```

2)Decision tree classification using weka

1)loading and preprocessing the data

Features:- outlook,temperature,humidity,windy

Label:-play



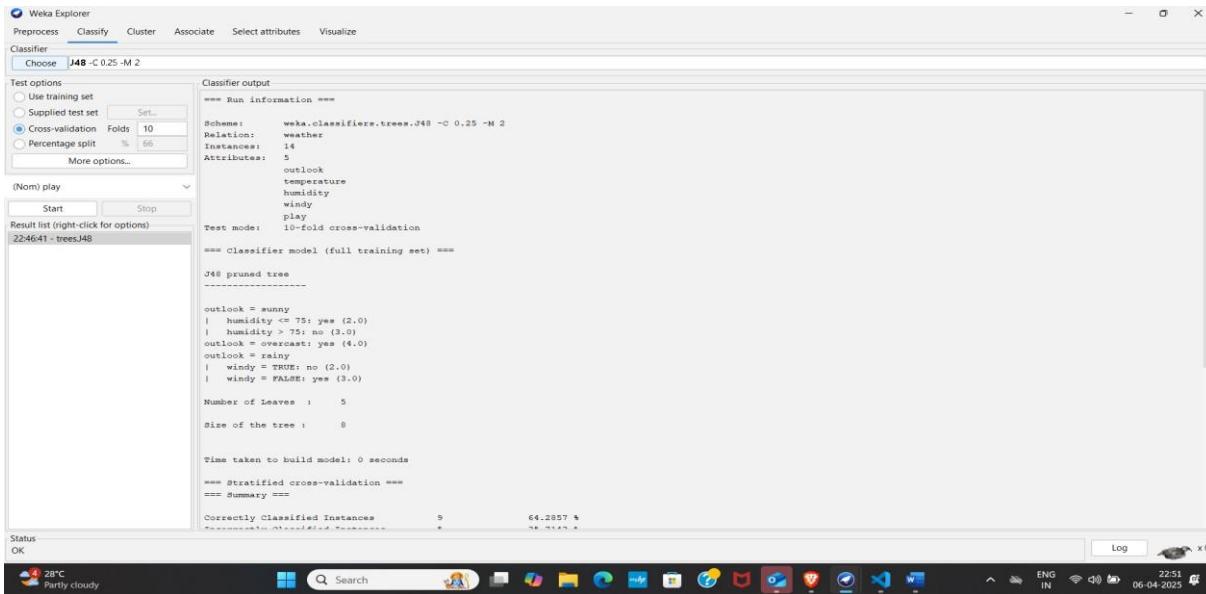
2) using J-48 decision tree

-> Click on classify in the tab above.

->Choose the classifier->trees->J48

->choose the values of cross-validation and percentage split

->Click on start.



```

Correctly Classified Instances      9          64.2857 %
Incorrectly Classified Instances   5          35.7143 %
Kappa statistic                   0.186
Mean absolute error               0.2857
Root mean squared error           0.4818
Relative absolute error            60         %
Root relative squared error       97.6586 %
Total Number of Instances         14

==== Detailed Accuracy By Class ====

      TP Rate  FP Rate  Precision  Recall   F-Measure  MCC     ROC Area  PRC Area  Class
0.778     0.600    0.700     0.778    0.737     0.189    0.789     0.847    yes
0.400     0.222    0.500     0.400    0.444     0.189    0.789     0.738    no
Weighted Avg.  0.643    0.465    0.629     0.643    0.632     0.189    0.789     0.808

==== Confusion Matrix ====

 a b  <-- classified as
7 2 | a = yes
3 2 | b = no

```

Result:-Accuracy of decision tree classification is 64.2857 %.

Advantages of decision tree classification-:

1. Easy to understand – Intuitive and visual.
2. No need for feature scaling – Works on raw data.
3. Handles both categorical & numerical data.
4. Captures non-linear patterns.
5. Feature importance – Highlights key predictors.
6. Fast prediction – Quick once trained.
7. Works for classification & regression.
8. No domain knowledge needed.

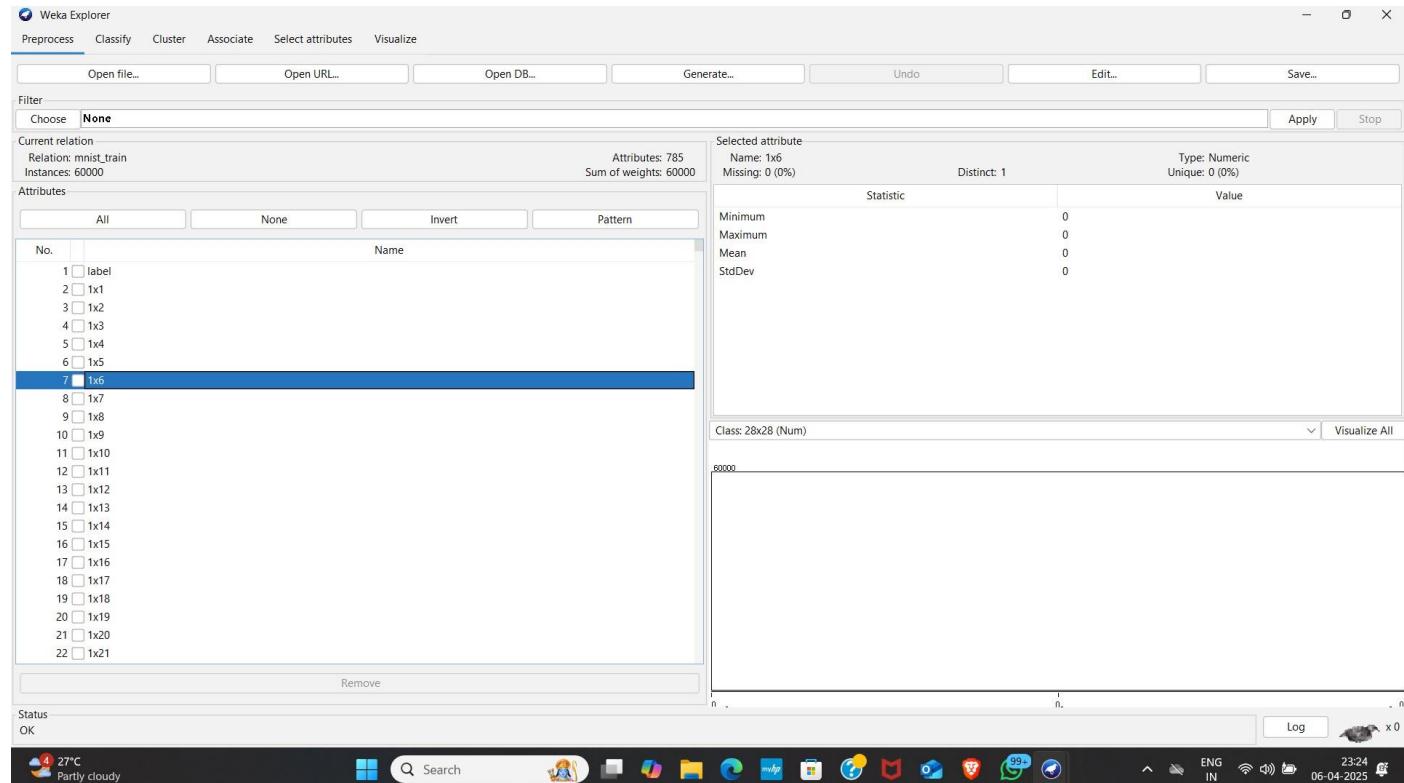
Disadvantages of decision tree classification-:

1. Overfitting – Trees can become too complex and fit noise.
2. Unstable – Small changes in data can change the structure.
3. Biased splits – Can favor features with more levels/categories.
4. Less accurate alone – Often outperformed by ensemble methods (like Random Forest).
5. Harder with continuous variables – Needs careful handling.

Assignment-5

1)K-Nearest neighbours classifiers for a dataset after pre-processing using weka/python libraries.

1)Loading the mnist_train.csv dataset into weka and doing exploratory analysis on the dataset. Features->786 pixels Labels->10(0-9).



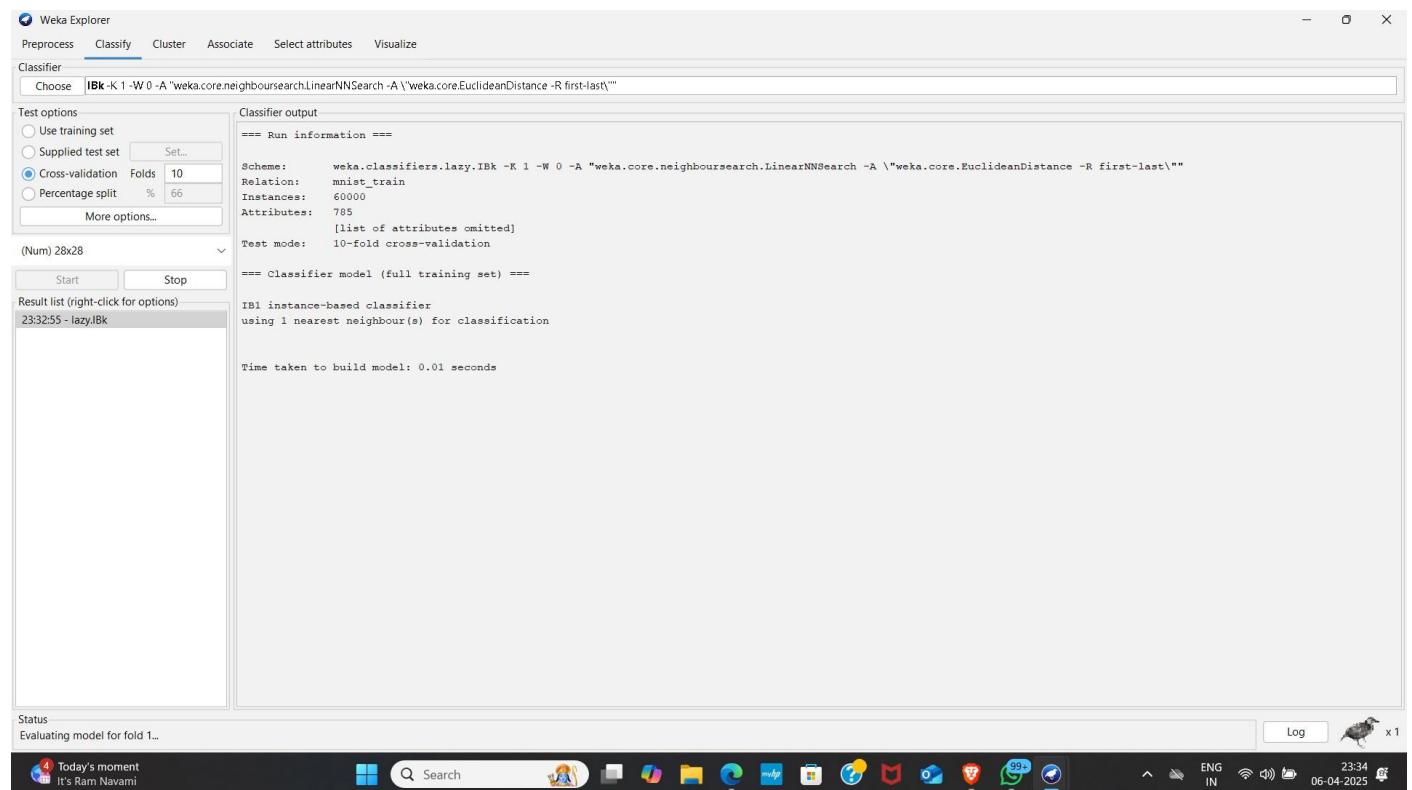
2)Normalize data

- >Go to the filters section again .
- >Choose the filter :Unsupervised->Attribute->Normalize.

3)Appying knn-classifier

- >Click on classifier in the tab.
- >choose the classifier:lazy->IBK.

Result:-



2) Python code for implementing k-nearest neighbours.

```
knns > 2.py > ...
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 # print(plt.style.available)
6 plt.style.use('seaborn-v0_8')
7 dfx=pd.read_csv('xdata.csv')
8 dfy=pd.read_csv('ydata.csv')
9
10 X=dfx.values
11 Y=dfy.values
12
13
14 X=X[:,1:]
15 Y=Y[:,1]
16
17
18
19 query_x=np.array([3,1])
20 plt.scatter(X[:,0],X[:,1],c=Y) # here we can also give c as array like like Y
21 plt.scatter(query_x[0],query_x[1],c='red')
22 plt.show()
23
24 def dist(x1,x2):
25     return np.sqrt(sum((x1-x2)**2))
26
27 def knn(X,Y,queryPoint,k=5):
28     vals=[]
29
30     m=X.shape[0]
31     for i in range(m):
32         d=dist(queryPoint,X[i])
33         vals.append((d,Y[i]))
34
35     vals=sorted(vals)
36
37     vals=vals[:k]
38     vals=np.array(vals)
39
40     # print(vals)
41
42     new_vals=np.unique(vals[:,1],return_counts=True)
43     print(new_vals)
44     index=new_vals[1].argmax()
45
46     pred=new_vals[0][index]
47
48     return pred
49
50
51 print(knn(X,Y,query_x))
52
53 # task to predict the accuracy of a knn
54
55 # Function to compute accuracy
56 def km_accuracy(X, Y, k=5, test_ratio=0.2):
57     # Flatten Y if needed
58     if Y.ndim > 1:
59         Y = Y.flatten()
60
61     # Shuffle data
62     indices = np.arange(X.shape[0])
63     np.random.shuffle(indices)
64     X = X[indices]
65
```

File Edit Selection View Go Run Terminal Help ← → ⌘ AI ML

EXPLORER .py 1.py recommended-system creating-recommending.py interview-question.py part-3.png formula.png 1.py knn 2.py knn

AI ML .pytest_cache data-aquisition-web_crawler-scr... data-visualization face-recognition knn knn-diabetes linear-algebra linear-regression myproject new_virtual_envi pandas probability-distribution-statistics... python language python language part2 python-coding-block recommended-system scrapyproject sklearn web_scraping web_automation-selenium 1.py decision-tree.py Diabetes_XTrain.csv Diabetes_YTrain.csv filename gender_submission.csv haarcascade_frontalface_alt.xml marks.csv mnist_train.csv movie_metadata.csv

OUTLINE TIMELINE KNN Accuracy: 100.0%

```

57 def knn_accuracy(X, Y, k=5, test_ratio=0.2):
58     X = X[indices]
59     Y = Y[indices]
60
61     # Split into train and test
62     split = int(len(X) * (1 - test_ratio))
63     X_train, Y_train = X[:split], Y[:split]
64     X_test, Y_test = X[split:], Y[split:]
65
66     correct = 0
67     for i in range(len(X_test)):
68         pred = knn(X_train, Y_train, X_test[i], k)
69         if pred == Y_test[i]:
70             correct += 1
71
72     acc = correct / len(X_test)
73     return acc
74
75 accuracy = knn_accuracy(X, Y, k=5, test_ratio=0.2)
76 print(f"KNN Accuracy: {round(accuracy * 100, 2)}%")
77
78
79
80
81
82
83

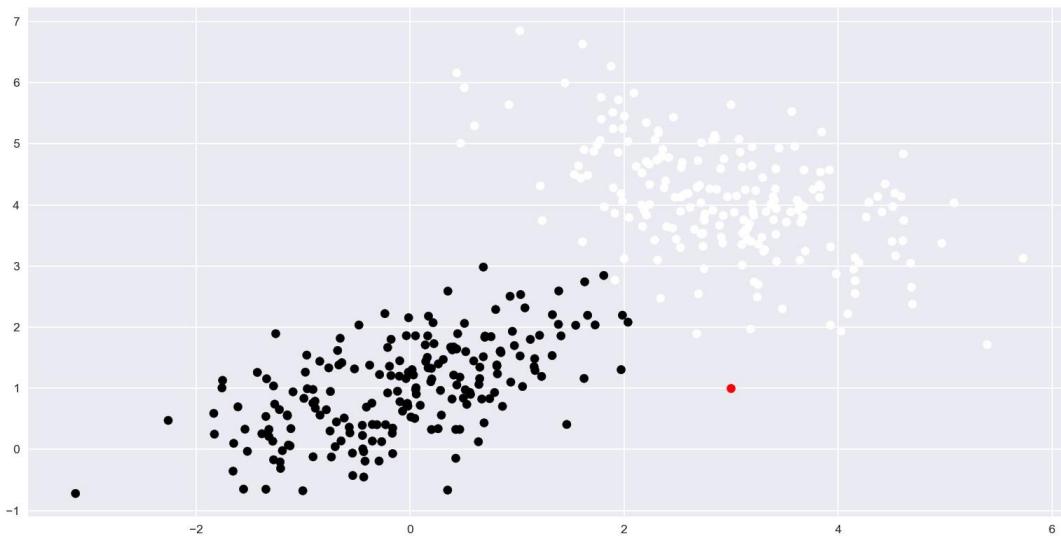
```

PROBLEMS OUTPUT TERMINAL PORTS

Ln 83, Col 52 Spaces: 4 UTF-8 CRLF {} Python 3.13.1 ENG IN 23:16 06-04-2025

Plot

Figure 1



Result:-As the dataset is clean and well separated we got an accuracy of 100 %

```
(array([0.]), array([5]))  
(array([0.]), array([5]))  
KNN Accuracy: 100.0%  
PS D:\AI ML> []
```

Applications of KNN Classifier (for Classification)

1. Medical Diagnosis

- Classifying diseases based on symptoms or test results.
- Example: Predicting if a tumor is **benign** or **malignant**.

2. Handwriting Recognition

- Classifying handwritten digits or letters (like in **MNIST dataset**).

3 Recommendation Systems

- Suggesting products based on user similarity (user-based KNN).
- Example: "**People who bought this also bought...**"

4. Image Classification

- Classifying images based on pixel similarity.
- Simple KNN can work well on small, clean datasets.

5. Spam Detection

- Classifying emails as spam or not based on content features.

6. Music Genre Classification

- Classify songs into genres based on audio features.

7. Customer Segmentation

- Grouping users by behavior or preferences for targeted marketing.

Assignment-6

1)K-mean Clustering for a dataset after pre-processing using weka/Python libraries.

Iris Dataset Summary

- **Rows:** 150 flower samples
- **Classes:** 3 types of Iris
 - *Setosa, Versicolor, Virginica*
- **Features:**
 - SepalLengthCm
 - SepalWidthCm
 - PetalLengthCm
 - PetalWidthCm
- **Label Column:** Species (used for comparison only in clustering)
- **Use Case:** Classification & Clustering (e.g., with K-Means)

Python code for K-mean clustering for iris data set using python

```
import pandas as pd

import numpy as np
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler

df = pd.read_csv('Iris.csv')
X = df.drop(columns=['Id', 'Species'])

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

X = np.array(X_scaled)
```

```

def kmeans(X, k, max_iters=100):
    np.random.seed(42)
    n_samples, n_features = X.shape

    centroids = X[np.random.choice(n_samples, k, replace=False)]

    for _ in range(max_iters):

        distances = np.linalg.norm(X[:, np.newaxis] - centroids, axis=2)
        cluster_labels = np.argmin(distances, axis=1)

        # Step 3: Recompute centroids
        new_centroids = np.array([X[cluster_labels == i].mean(axis=0) for i in range(k)])

        if np.allclose(centroids, new_centroids):
            break

        centroids = new_centroids

    return cluster_labels, centroids

# Run K-Means with k=3 (for Iris dataset)
# Run K-Means
k = 3
labels, centroids = kmeans(X, k)

# Add cluster labels to the DataFrame
df['Cluster'] = labels

print(df[['SepalLengthCm', 'SepalWidthCm', 'Cluster']].head())

```

```

# Optional: visualize clusters (based on 1st 2 features)

plt.figure(figsize=(8, 6))

plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis', s=50)

plt.scatter(centroids[:, 0], centroids[:, 1], color='red', marker='X', s=200, label='Centroids')

plt.title('K-Means Clustering on Iris Dataset')

plt.xlabel('Sepal Length (standardized)')

plt.ylabel('Sepal Width (standardized)')

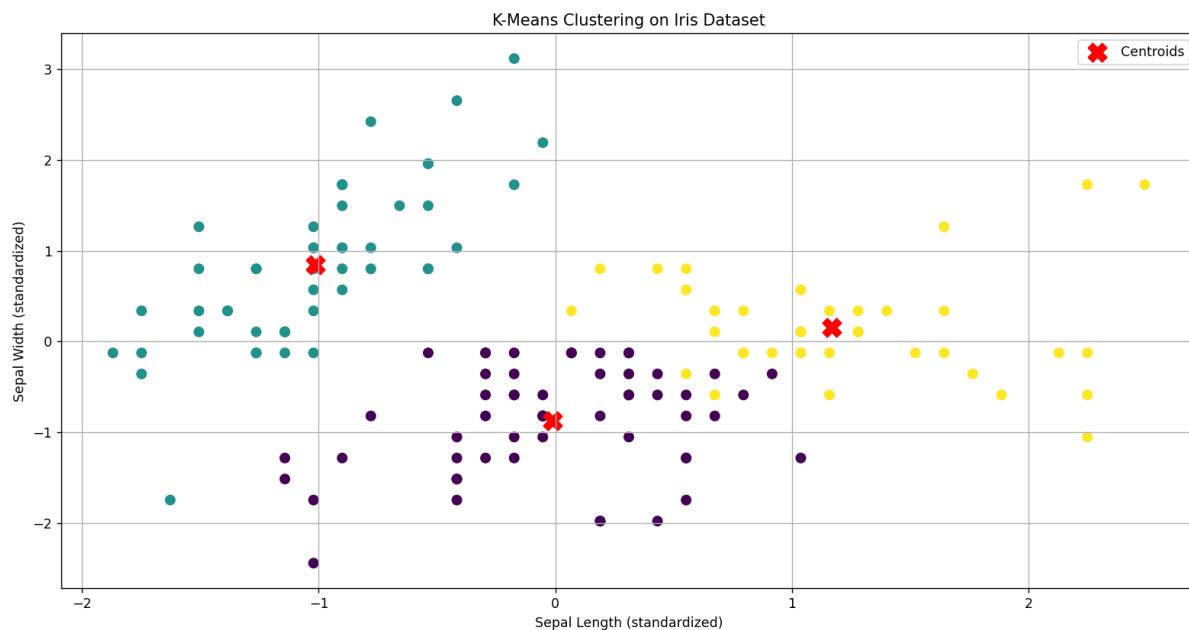
plt.legend()

plt.grid(True)

plt.show()

```

Plot



Result:-

	SepalLengthCm	SepalWidthCm	Cluster
0	5.1	3.5	1
1	4.9	3.0	1
2	4.7	3.2	1
3	4.6	3.1	1
4	5.0	3.6	1

Top Applications of K-Means Clustering

1. Customer Segmentation

- **Goal:** Group customers based on purchasing behavior, demographics, or preferences.
- **Example:** E-commerce platforms use K-Means to segment users for targeted marketing.

2. Market Basket Analysis / Product Recommendation

- **Goal:** Group products that are often bought together or group users with similar purchase patterns.
- **Example:** Amazon or Netflix uses it to recommend products/movies.

3. Document or Text Clustering

- **Goal:** Automatically group similar documents or articles.
- **Example:** News aggregation, topic modeling, grouping tweets, etc.

4. Image Compression / Segmentation

- **Goal:** Reduce colors in an image by clustering similar color pixels.
- **Example:** Use K-Means to replace similar colors with their centroid → smaller file size with preserved visual quality.

5. Education Analytics

- **Goal:** Group students based on performance, engagement, or learning styles.

Assignment-7

Title: Agglomerative (Bottom-Up) Clustering for a Dataset After Preprocessing using Python Libraries

Objective: To implement and evaluate Agglomerative (bottom-up) clustering on a preprocessed dataset using Python libraries. The experiment involves data pre-processing followed by the application of Agglomerative clustering to identify natural groupings within the dataset.

Dataset: new.csv

Description: The dataset relates to marketing campaigns of a Portuguese banking institution, consisting of demographic and contact attributes along with campaign outcomes.

Tools used:- Python 3.x

Libraries used:-

Pandas, numpy , matplotlib, seaborn, sklearn (StandardScaler, LabelEncoder, AgglomerativeClustering) , scipy (linkage, dendrogram) .

Code

```
# Importing libraries import pandas as pd import numpy as np  
import matplotlib.pyplot as plt import seaborn as sns  
from sklearn.preprocessing import StandardScaler, LabelEncoder from sklearn.cluster import  
AgglomerativeClustering from scipy.cluster.hierarchy import dendrogram, linkage
```

Load the dataset

```
# Importing libraries import pandas as pd import numpy as np  
import matplotlib.pyplot as plt import seaborn as sns  
from sklearn.preprocessing import StandardScaler, LabelEncoder from sklearn.cluster import  
AgglomerativeClustering from scipy.cluster.hierarchy import dendrogram, linkage
```

Load the dataset

```
df = pd.read_csv( '/mnt/data/bank-full - bank-full.csv.csv' , sep= ';' )
```

Tools Used:

Procedure:

```
# Display basic information print (df.info()) print (df.head())
```

```
# Pre-processing
```

```
# Encode categorical variables
```

```
label_encoders = {} for column in df.select_dtypes(include=[ 'object' ]).columns:
```

```
le = LabelEncoder() df[column] = le.fit_transform(df[column]) label_encoders[column] = le
```

```
# Feature scaling scaler = StandardScaler() df_scaled = scaler.fit_transform(df)
```

```
# Create Linkage matrix for dendrogram linked = linkage(df_scaled, method= 'ward' )
```

```
# Plot dendrogram
```

```
plt.figure(figsize=( 12 , 6 )) dendrogram(linked, truncate_mode= 'level' , p= 5 ) plt.title( 'Dendrogram for Hierarchical Clustering' ) plt.xlabel( 'Data points' ) plt.ylabel( 'Euclidean distance' ) plt.show()
```

```
# Apply Agglomerative Clustering agg_cluster = AgglomerativeClustering(n_clusters= 3 , affinity= 'euclidean' , linkage= 'ward' ) cluster_labels = agg_cluster.fit_predict(df_scaled)
```

```
# Add cluster labels to the dataframe df[ 'cluster' ] = cluster_labels
```

```
# Analyze Results print ( f "Unique clusters found: {np.unique(cluster_labels)} " ) print (df[ 'cluster' ].value_counts())
```

```
# Visualization of clusters  
plt.figure(figsize=( 10 , 6 ))  
  
sns.scatterplot( x=df_scaled[:, 0 ],  
y=df_scaled[:, 1 ], hue=cluster_labels, palette= 'Set2' , legend= 'full'  
)  
  
plt.title( 'Agglomerative Clustering Visualization' ) plt.xlabel( 'Feature 1 (Standardized)' ) plt.ylabel(  
'Feature 2 (Standardized)' )  
  
plt.show()
```

The image shows two screenshots of the Google Colab interface, one above the other, demonstrating a Python script for hierarchical clustering on the 'bank-full' dataset.

Code Cell 1 (Top):

```

# Importing libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.cluster import AgglomerativeClustering
from scipy.cluster.hierarchy import dendrogram, linkage

# Load the dataset
df=pd.read_csv("/content/bank-full - bank-full.csv.csv")

# Display basic information
print(df.info())
print(df.head())

# Pre-processing
# Encode categorical variables
label_encoders = {}
for column in df.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    df[column] = le.fit_transform(df[column])
    label_encoders[column] = le

# Feature scaling
scaler = StandardScaler()
df_scaled = scaler.fit_transform(df)

# Create linkage matrix for dendrogram
linked = linkage(df_scaled, method='ward')

```

A tooltip at the bottom left of the code cell says "Your session crashed after using all available RAM." and "View runtime logs".

Code Cell 2 (Bottom):

```

df_scaled = scaler.fit_transform(df)

# Create Linkage matrix for dendrogram
linked = linkage(df_scaled, method='ward')

# Plot dendrogram
plt.figure(figsize=(12, 6))
dendrogram(linked, truncate_mode='level', p=5)
plt.title("Dendrogram for Hierarchical Clustering")
plt.xlabel('Data points')
plt.ylabel('Euclidean distance')
plt.show()

# Apply Agglomerative Clustering
agg_cluster = AgglomerativeClustering(n_clusters=3, affinity='euclidean', linkage='ward')
cluster_labels = agg_cluster.fit_predict(df_scaled)

# Add cluster labels to the dataframe
df['cluster'] = cluster_labels

# Analyze Results
print(f"Unique clusters found: {np.unique(cluster_labels)}")
print(df['cluster'].value_counts())

# visualization of clusters
plt.figure(figsize=(10, 6))
sns.scatterplot(
    x=df_scaled[:, 0],
    y=df_scaled[:, 1],
    hue=cluster_labels,
    palette='Set2'
)

```

A tooltip at the bottom left of the code cell says "Your session crashed after using all available RAM." and "View runtime logs".

```

# Visualization of clusters
plt.figure(figsize=(10, 6))
sns.scatterplot(
    x=df_scaled[:, 0],
    y=df_scaled[:, 1],
    hue=cluster_labels,
    palette='Set2',
    legend='full'
)
plt.title('Agglomerative Clustering Visualization')
plt.xlabel('Feature 1 (Standardized)')
plt.ylabel('Feature 2 (Standardized)')
plt.show()

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 45211 entries, 0 to 45210
Data columns (total 17 columns):
 #   Column      Non-Null Count  Dtype  
 --- 
 0   age         45211 non-null   int64  
 1   job          45211 non-null   object  
 2   marital     45211 non-null   object  
 3   education   45211 non-null   object  
 4   default     45211 non-null   object  
 5   balance     45211 non-null   int64  
 6   housing     45211 non-null   object  
 7   loan         45211 non-null   object  
 8   contact     45211 non-null   object  
 9   day          45211 non-null   int64  
 10  month        45211 non-null   object  
 11  duration    45211 non-null   int64  
 12  campaign    45211 non-null   int64  
 13  pdays       45211 non-null   int64  
 14  previous    45211 non-null   int64  
 15  poutcome    45211 non-null   object  
 16  y           45211 non-null   object  
dtypes: int64(7), object(10)
memory usage: 5.9+ MB
None
   age      job marital education default balance housing loan \
0  58  management  married   tertiary    no   2143   yes   no
1  44  technician single   secondary   no     29   yes   no
2  33  entrepreneur married   secondary   no      2   yes   yes
3  47  blue-collar married unknown    no   1506   yes   no
4  33      unknown single   unknown    no      1   no   no

   contact day month duration campaign pdays previous poutcome y
0  unknown   5   may    261        1   -1      0  unknown   no
1  unknown   5   may    151        1   -1      0  unknown   no
2  unknown   5   may    76         1   -1      0  unknown   no
3  unknown   5   may    92         1   -1      0  unknown   no
4  unknown   5   may   198        1   -1      0  unknown   no

```

Your session crashed after using all available RAM. View runtime logs

Assignment-8

Objective

To implement and evaluate DBSCAN clustering on a pre-processed dataset using Python libraries. The experiment involves data pre-processing steps followed by the application of DBSCAN to identify clusters and noise

Dataset Description

Name: bank-full.csv

Description: Dataset related to direct marketing campaigns of a Portuguese banking institution. It includes attributes such as age, job, marital status, education, default status, housing loan, and personal loan.

Code

```
# Load the dataset  
df = pd.read_csv( '/mnt/data/bank-full - bank-full.csv.csv' , sep= ';' )  
  
# Display basic information  
print (df.info())  
  
# Display first few rows  
print (df.head())
```

Procedure

Step 1: Import Required Libraries

```
import pandas as pd
```

```
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
from sklearn.preprocessing import StandardScaler, LabelEncoder  
from sklearn.cluster import DBSCAN
```

Step 2: Load and Explore Dataset

Step 3: Pre-processing the Data

Encoding categorical variables using Label Encoding.

Scaling the features using StandardScaler

```
# Encode categorical variables  
  
label_encoders = {}  
  
for column in df.select_dtypes(include=[ 'object' ]).columns:  
    le = LabelEncoder()  
    df[column] = le.fit_transform(df[column])  
    label_encoders[column] = le  
  
  
# Feature scaling  
  
scaler = StandardScaler()  
df_scaled = scaler.fit_transform(df)
```

Step 4: Apply DBSCAN Clustering

```
# Applying DBSCAN algorithm  
  
dbscan = DBSCAN(eps= 2 , min_samples= 5 )  
dbscan.fit(df_scaled)
```

```

# Assigning cluster labels
df[ 'cluster' ] = dbscan.labels_

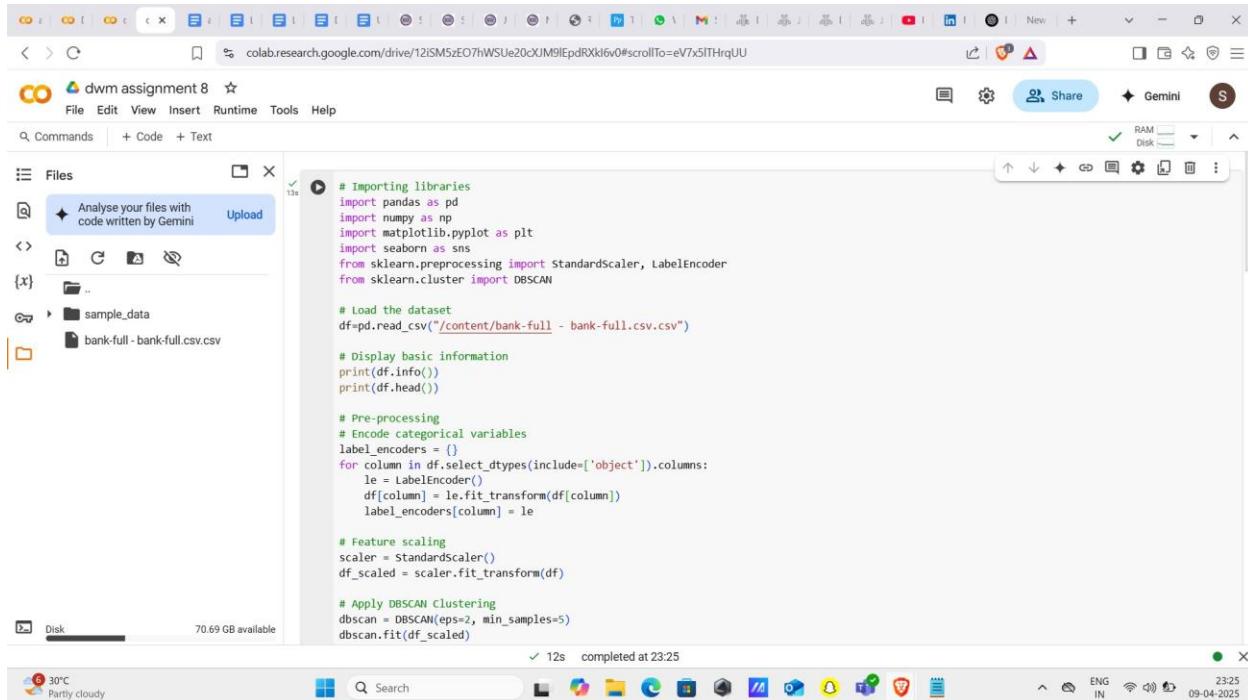

Step 5: Analyze and Visualize the Results

# Unique clusters
print ( f "Unique clusters found: {np.unique(dbscan.labels_)}" )

# Count of points in each cluster
print (df[ 'cluster' ].value_counts())


# Visualization of clusters
plt.figure(figsize=( 10 , 6 ))
sns.scatterplot(
    x=df_scaled[:, 0 ],
    y=df_scaled[:, 1 ],
    hue=dbscan.labels_,
    palette= 'Set1',
    legend= 'full'
)
plt.title( 'DBSCAN Clustering Visualization' )
plt.xlabel( 'Feature 1 (Standardized)' )
plt.ylabel( 'Feature 2 (Standardized)' )
plt.show()

```



The screenshot shows a Google Colab notebook titled "dwm assignment 8". The code cell contains Python code for DBSCAN clustering. The code imports libraries, reads a dataset, encodes categorical variables, scales features, applies DBSCAN, adds cluster labels, and visualizes the results. The status bar at the bottom indicates the code completed at 23:25.

```

# Importing libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.cluster import DBSCAN

# Load the dataset
df=pd.read_csv("/content/bank-full - bank-full.csv.csv")

# Display basic information
print(df.info())
print(df.head())

# Pre-processing
# Encode categorical variables
label_encoders = {}
for column in df.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    df[column] = le.fit_transform(df[column])
    label_encoders[column] = le

# Feature scaling
scaler = StandardScaler()
df_scaled = scaler.fit_transform(df)

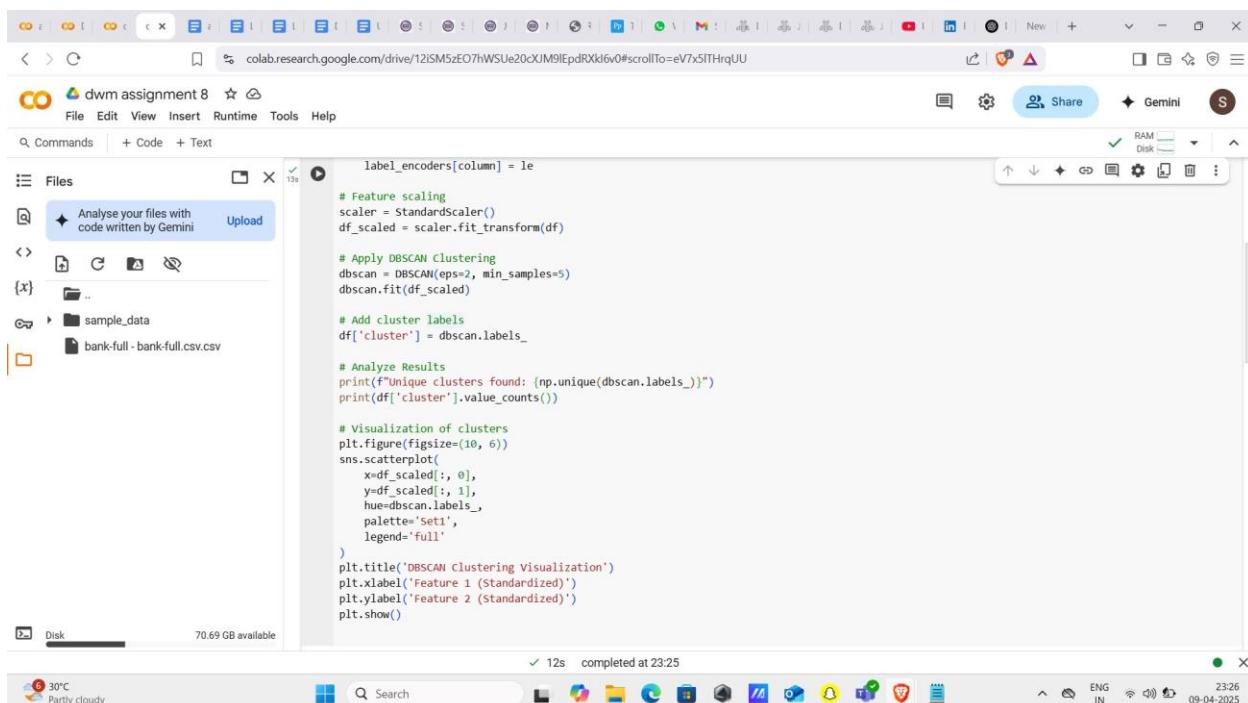
# Apply DBSCAN Clustering
dbscan = DBSCAN(eps=2, min_samples=5)
dbscan.fit(df_scaled)

# Add cluster labels
df['cluster'] = dbscan.labels_

# Analyze Results
print(f"Unique clusters found: {np.unique(dbscan.labels_)}")
print(df['cluster'].value_counts())

# Visualization of clusters
plt.figure(figsize=(10, 6))
sns.scatterplot(
    x=df_scaled[:, 0],
    y=df_scaled[:, 1],
    hue=dbscan.labels_,
    palette='Set1',
    legend='full'
)
plt.title('DBSCAN Clustering Visualization')
plt.xlabel('Feature 1 (Standardized)')
plt.ylabel('Feature 2 (Standardized)')
plt.show()

```



The screenshot shows the same Google Colab notebook. The code cell now includes the final visualization command. The status bar at the bottom indicates the code completed at 23:26.

dwm assignment 8

File Edit View Insert Runtime Tools Help

Commands + Code + Text

Files

Analyse your files with code written by Gemini

Upload

{x} sample_data bank-full - bank-full.csv.csv

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 45211 entries, 0 to 45210
Data columns (total 17 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   age         45211 non-null   int64  
 1   job          45211 non-null   object  
 2   marital     45211 non-null   object  
 3   education   45211 non-null   object  
 4   default     45211 non-null   object  
 5   balance     45211 non-null   int64  
 6   housing     45211 non-null   object  
 7   loan         45211 non-null   object  
 8   contact     45211 non-null   object  
 9   day          45211 non-null   int64  
 10  month        45211 non-null   object  
 11  duration    45211 non-null   int64  
 12  campaign    45211 non-null   int64  
 13  pdays       45211 non-null   int64  
 14  previous    45211 non-null   int64  
 15  poutcome    45211 non-null   object  
 16  y           45211 non-null   object  
dtypes: int64(7), object(10)
memory usage: 5.9+ MB
None
   age   job   marital   education   default   balance   housing   loan \
0  58   management   married   tertiary   no   2143   yes   no
1  44   technician   single   secondary   no    29   yes   no
2  33   entrepreneur   married   secondary   no     2   yes   yes
3  47   blue-collar   married   unknown   no   1506   yes   no
4  33   unknown     single   unknown   no     1   no   no
   contact   day   month   duration   campaign   pdays   previous   poutcome   y
0  unknown   5   may    261        1      -1      0   unknown   no
1  unknown   5   may    151        1      -1      0   unknown   no
2  unknown   5   may    76         1      -1      0   unknown   no
3  unknown   5   may    92         1      -1      0   unknown   no
4  unknown   5   may   198        1      -1      0   unknown   no
23 24 25 26 27 28 29 30 31 32 33 34 35]
cluster
0   18180
2   13068
1   3694
-1   3650
5   2422
3   2198
4   1089
7   206
19   116
9   106
17   84
18   71
23   66
6   46
16   29
8   28
32   18
27   18

```

Disk 70.69 GB available ✓ 12s completed at 23:25

30°C Partly cloudy

Search

ENG IN 23:26 09-04-2025

dwm assignment 8

File Edit View Insert Runtime Tools Help

Commands + Code + Text

Files

Analyse your files with code written by Gemini

Upload

{x} sample_data bank-full - bank-full.csv.csv

```
None
   age   job   marital   education   default   balance   housing   loan \
0  58   management   married   tertiary   no   2143   yes   no
1  44   technician   single   secondary   no    29   yes   no
2  33   entrepreneur   married   secondary   no     2   yes   yes
3  47   blue-collar   married   unknown   no   1506   yes   no
4  33   unknown     single   unknown   no     1   no   no
   contact   day   month   duration   campaign   pdays   previous   poutcome   y
0  unknown   5   may    261        1      -1      0   unknown   no
1  unknown   5   may    151        1      -1      0   unknown   no
2  unknown   5   may    76         1      -1      0   unknown   no
3  unknown   5   may    92         1      -1      0   unknown   no
4  unknown   5   may   198        1      -1      0   unknown   no
23 24 25 26 27 28 29 30 31 32 33 34 35]
cluster
0   18180
2   13068
1   3694
-1   3650
5   2422
3   2198
4   1089
7   206
19   116
9   106
17   84
18   71
23   66
6   46
16   29
8   28
32   18
27   18

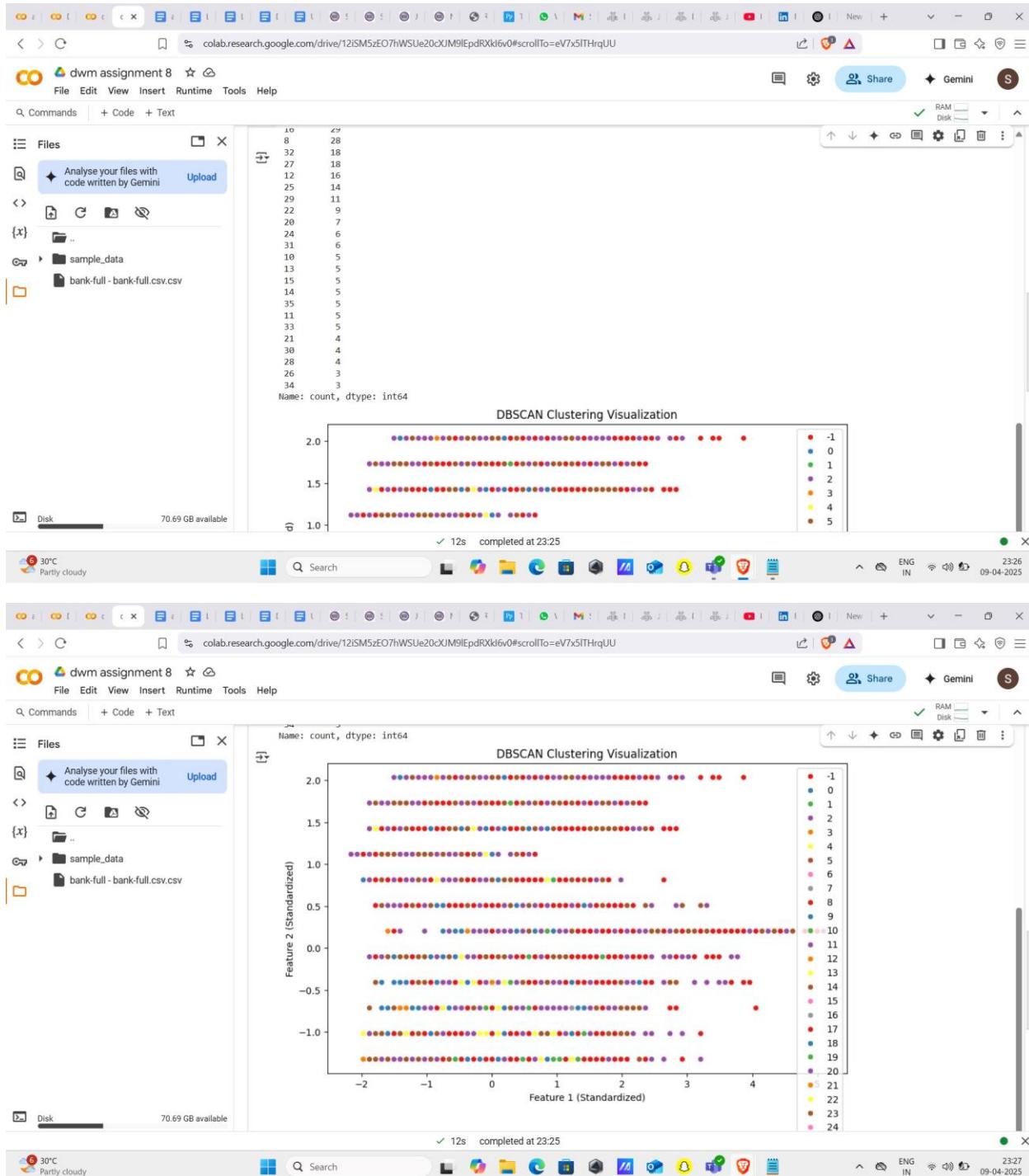
```

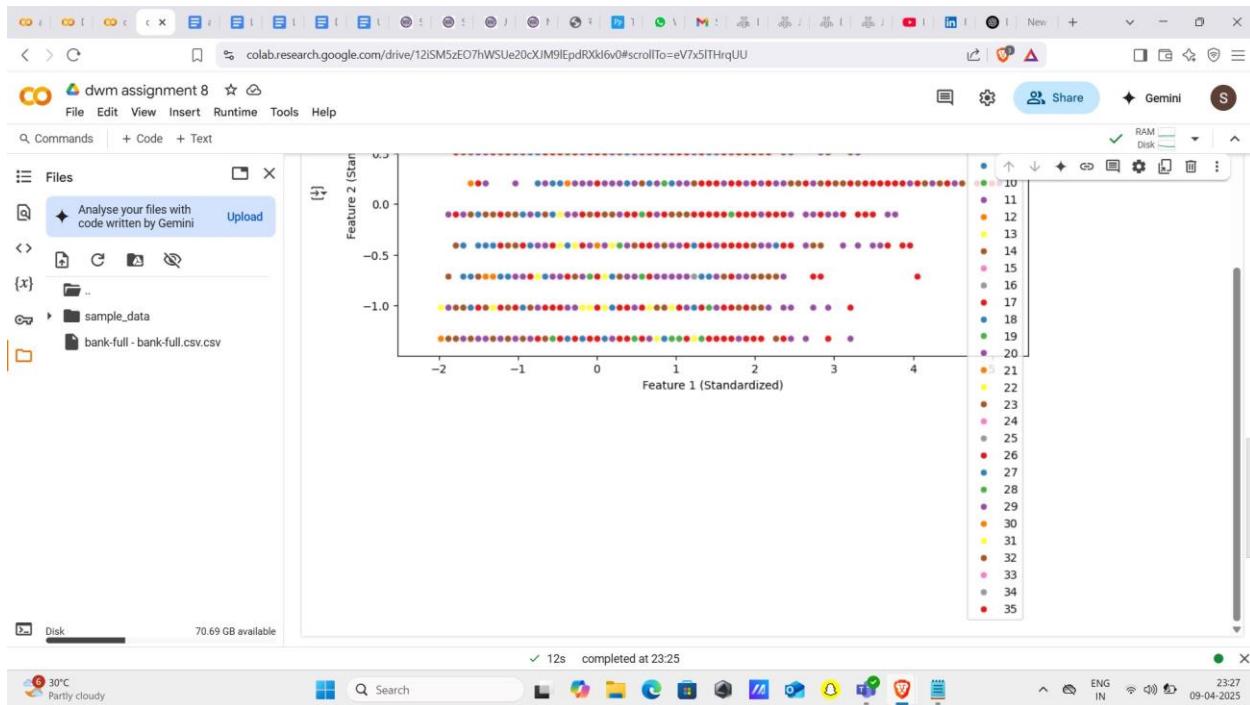
Disk 70.69 GB available ✓ 12s completed at 23:25

30°C Partly cloudy

Search

ENG IN 23:26 09-04-2025





Observations

DBSCAN divided the data into clusters and identified noise points.

Cluster -1 corresponds to noise or outliers.

DBSCAN does not require specifying the number of clusters beforehand.

Clustering results are sensitive to the parameters `eps` (neighborhood size) and `min_samples`

Assignment-9

1)Apriori algorithm for association rule mining using Python libraries.

Dataset Overview

- Format: CSV (Comma-Separated Values)
- Rows: 1,4965 (based on the visible data)
- Columns: 11 (labeled 0 through 10)
- Data Type: Categorical (item names as strings)
- Purpose: Likely used for market basket analysis, such as identifying purchasing patterns, frequent itemsets, or association rules (e.g., "if a customer buys X, they are likely to buy Y").

Python code for Apriori algorithm for association rule mining.

```
import pandas as pd
from mlxtend.frequent_patterns import apriori, association_rules

# Step 1: Load the dataset
try:
    data = pd.read_csv('basket.csv', header=None)
    print(f"Dataset loaded successfully. Shape: {data.shape}")
except FileNotFoundError:
    print("Error: 'basket.csv' not found. Please ensure the file is in the correct directory.")
    exit(1)

# Step 2: Preprocess the data
# Convert the dataset into a list of transactions
transactions = []
for i in range(len(data)):
    # Extract non-null items from each row
    transaction = [item for item in data.iloc[i] if pd.notna(item)]
    transactions.append(transaction)

# Get all unique items
all_items = sorted(set(item for transaction in transactions for item in transaction))
print(f"Number of unique items: {len(all_items)}")
print(f"Number of transactions: {len(transactions)}")

# Create a one-hot encoded DataFrame
one_hot_data = pd.DataFrame(0, index=range(len(transactions)), columns=all_items)
for i, transaction in enumerate(transactions):
    for item in transaction:
        one_hot_data.at[i, item] = 1
```

```

# Print a sample of the one-hot encoded data
print("\nSample of one-hot encoded data (first 5 rows):")
print(one_hot_data.head())

# Step 3: Apply the Apriori algorithm
# Lower the minimum support to 1% (adjust as needed)
min_support = 0.01 # Changed from 0.75 to 0.01
frequent_itemsets = apriori(one_hot_data, min_support=min_support, use_colnames=True)

# Check if frequent itemsets were found
if frequent_itemsets.empty:
    print(f"\nNo frequent itemsets found with min_support={min_support}. Try lowering the support threshold.")
else:
    print("\nFrequent Itemsets:")
    print(frequent_itemsets)

# Step 4: Generate association rules (only if frequent itemsets exist)
min_confidence = 0.1
if not frequent_itemsets.empty:
    rules = association_rules(frequent_itemsets, metric="confidence",
                               min_threshold=min_confidence)

# Step 5: Evaluate and display the results
print("\nAssociation Rules:")
if rules.empty:
    print(f"No association rules found with min_confidence={min_confidence}.")
else:
    print(rules[['antecedents', 'consequents', 'support', 'confidence', 'lift']])

# Save results to CSV files
frequent_itemsets.to_csv('frequent_itemsets.csv', index=False)
rules.to_csv('association_rules.csv', index=False)

# Basic evaluation metrics
print("\nNumber of frequent itemsets:", len(frequent_itemsets))
print("Number of association rules:", len(rules))
print("\nTop 5 rules sorted by lift:")
print(rules[['antecedents', 'consequents', 'support', 'confidence', 'lift']].sort_values(by='lift',
                                                                                     ascending=False).head())
else:
    print("\nSkipping association rules generation due to empty frequent itemsets.")

```

Result:-

->With 75 % support and 85% confidence we can't get any rule

```
Dataset loaded successfully. Shape: (14964, 11)
Number of unique items: 178
Number of transactions: 14964

Sample of one-hot encoded data (first 5 rows):
   0 1 10 2 3 4 5 6 7 8 9 ... tropical fruit turkey vinegar waffles whipped/sour cream whisky white bread white wine whole milk yogurt zwieback
0  0 1 1 1 1 1 1 1 1 1 1 ...      0     0     0     0           0     0     0     0     0     0     0     0
1  0 0 0 0 0 0 0 0 0 0 0 ...      0     0     0     0           0     0     0     0     1     0     0
2  0 0 0 0 0 0 0 0 0 0 0 ...      0     0     0     0           0     0     0     0     1     1     0
3  0 0 0 0 0 0 0 0 0 0 0 ...      0     0     0     0           0     0     0     0     0     0     0
4  0 0 0 0 0 0 0 0 0 0 0 ...      0     0     0     0           0     0     0     0     0     0     0

[5 rows x 178 columns]
C:\Users\labhg\AppData\Roaming\Python\Python313\site-packages\mlxtend\frequent_patterns\fpcommon.py:161: DeprecationWarning: DataFrames with non-bool types result in worse computational performance and their support might be discontinued in the future. Please use a DataFrame with bool type
  warnings.warn(
No frequent itemsets found with min_support=0.75. Try lowering the support threshold.

Skipping association rules generation due to empty frequent itemsets.
PS D:\VAI ML>
```

->Result on lowering the support and confidence we get.

```
C:\Users\labhg\AppData\Roaming\Python\Python313\site-packages\mlxtend\frequent_patterns\fpcommon.py:161: DeprecationWarning: DataFrames with non-bool types result in worse computational performance and their support might be discontinued in the future. Please use a DataFrame with bool type
  warnings.warn(
Frequent Itemsets:
    support          itemsets
0  0.021385      (UHT-milk)
1  0.033948      (beef)
2  0.021786      (berries)
3  0.016573      (beverages)
4  0.045399      (bottled beer)
.. ...
64 0.010559  (other vegetables, rolls/buns)
65 0.014836  (other vegetables, whole milk)
66 0.013967  (rolls/buns, whole milk)
67 0.011628  (soda, whole milk)
68 0.011160  (yogurt, whole milk)

[69 rows x 2 columns]

Association Rules:
   antecedents  consequents  support  confidence  lift
0  (other vegetables)  (whole milk)  0.014836  0.121511  0.769482
1  (rolls/buns)        (whole milk)  0.013967  0.126974  0.804082
2  (soda)              (whole milk)  0.011628  0.119752  0.758346
3  (yogurt)            (whole milk)  0.011160  0.129961  0.822995
```

2) Apriori algorithm for association rule mining using weka.

The basket.csv dataset is in a transactional format (items listed in columns per transaction), but Weka's Apriori algorithm requires data in a specific format, typically **ARFF (Attribute-Relation File Format)** with binary attributes (1 for item presence, 0 for absence) or a sparse format. so we need to preprocess it.

Python code for converting basket.csv to basket.arff.

```
import pandas as pd
from collections import defaultdict

# Load the dataset
df = pd.read_csv('basket.csv', header=None)

# Get all unique items
items = df.values.flatten()
unique_items = pd.Series(items).dropna().unique()

# Create a binary matrix (1 if item is present, 0 otherwise)
binary_data = pd.DataFrame(0, index=df.index, columns=unique_items)
for idx, row in df.iterrows():
    for item in row.dropna():
        binary_data.loc[idx, item] = 1

# Save to CSV (intermediate step)
binary_data.to_csv('basket_binary.csv', index=False)

# Convert to ARFF
arff_content = f"""\@RELATION market_basket
"""

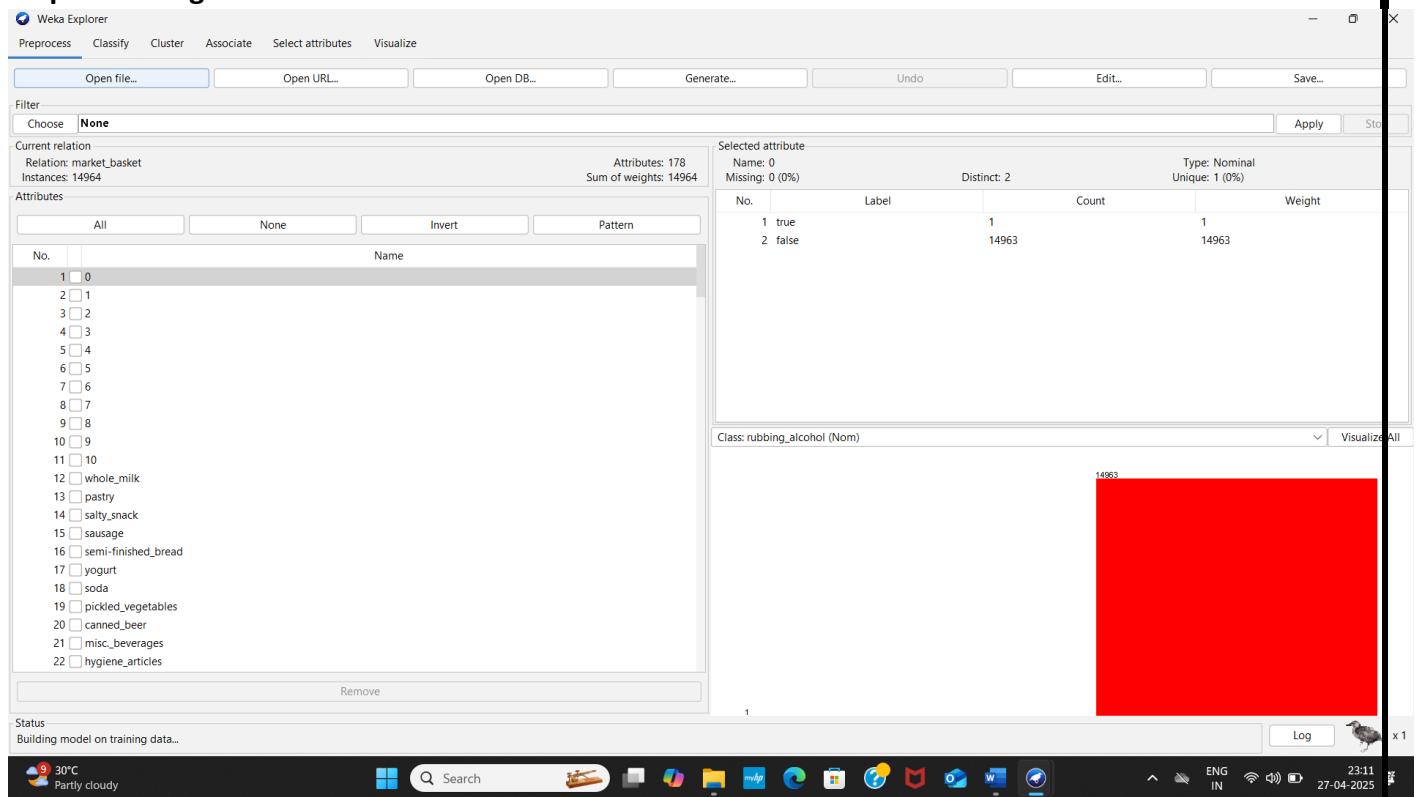
for item in unique_items:
    arff_content += f"\@ATTRIBUTE {item.replace(' ', '_')} {{true, false}}\n"

arff_content += "\n@\n@DATA\n"
for _, row in binary_data.iterrows():
    row_values = ['true' if val == 1 else 'false' for val in row]
    arff_content += ','.join(row_values) + '\n'

# Write to ARFF file
with open('basket.arff', 'w') as f:
    f.write(arff_content)

print("ARFF file 'basket.arff' created.")
```

Step-2 Loading the dataset into weka



Step-3 Switch to the Associate Tab:

- In Weka Explorer, click the “Associate” tab.

Select Apriori:

- Click the “Choose” button and select weka.associations.Apriori from the list.

Configure Apriori Parameters:

- Click the “Apriori” text to open the configuration window.
- Key parameters to adjust:
 - minSupport:** Minimum support threshold (e.g., 0.05 means an itemset must appear in 5% of transactions, or ~75 transactions). Start with a low value (e.g., 0.01–0.1) to capture frequent itemsets.
 - minMetric:** Minimum confidence for rules (e.g., 0.9 means 90% confidence that the rule holds). Common values are 0.5–0.9.
 - numRules:** Maximum number of rules to generate (e.g., 10–100). Set to 10 for initial testing.
 - metricType:** Choose “Confidence” (default) for standard association rules.

Step-4 Run the algorithm.

Click “Start” to execute Apriori.

Result:-

Weka Explorer

Preprocess Classify Cluster Associate **Select attributes** Visualize

Associate **Choose Apriori -N 10 -T 0 -C 0.85 -D 0.05 -U 1.0 -M 0.75 -S 1.0 -c -1**

Start Stop

Result list (right-click for ...)

23:09:08 - Apriori
23:29:17 - Apriori
23:36:43 - Apriori

Associate output

```
== Run information ==
Scheme: weka_associations.Apriori -N 10 -T 0 -C 0.85 -D 0.05 -U 1.0 -M 0.75 -S 1.0 -c -1
Relation: market_basket-weka.filters.unsupervised.attribute.Remove-R1-12,21-178
Instances: 14964
Attributes: 8
pastry
salty_snack
sausage
semi-finished_bread
yogurt
soda
pickled_vegetables
canned_beer
== Associate model (full training set) ==

Apriori
=====
Minimum support: 0.95 (14216 instances)
Minimum metric <confidence>: 0.85
Number of cycles performed: 1

Generated sets of large itemsets:
Size of set of large itemsets L(1): 4
Size of set of large itemsets L(2): 3
Size of set of large itemsets L(3): 1

Best rules found:

1. salty_snack=false 14683 => pickled_vegetables=false 14552 <conf:(0.99)> lift:(1) lev:(0) [0] conv:(1)
2. salty_snack=false semi-finished_bread=false 14546 => pickled_vegetables=false 14415 <conf:(0.99)> lift:(1) lev:(-0) [0] conv:(0.99)
3. semi-finished_bread=false 14822 => pickled_vegetables=false 14688 <conf:(0.99)> lift:(1) lev:(-0) [-1] conv:(0.98)
4. salty_snack=false 14683 => semi-finished_bread=false 14546 <conf:(0.99)> lift:(1) lev:(0) [2] conv:(1.01)
5. salty_snack=false pickled_vegetables=false 14552 => semi-finished_bread=false 14415 <conf:(0.99)> lift:(1) lev:(0) [1] conv:(1)
6. pickled_vegetables=false 14830 => semi-finished_bread=false 14688 <conf:(0.99)> lift:(1) lev:(-0) [-1] conv:(0.98)
7. salty_snack=false 14683 => semi-finished_bread=false pickled_vegetables=false 14415 <conf:(0.99)> lift:(1) lev:(0) [2] conv:(1.01)
8. semi-finished_bread=false pickled_vegetables=false 14688 => salty_snack=false 14415 <conf:(0.99)> lift:(1) lev:(0) [2] conv:(1.01)
9. semi-finished_bread=false 14822 => salty_snack=false 14546 <conf:(0.98)> lift:(1) lev:(0) [2] conv:(1)
10. pickled_vegetables=false 14830 => salty_snack=false 14552 <conf:(0.98)> lift:(1) lev:(0) [0] conv:(1)
```

Status OK Log

29°C Partly cloudy Search 23:37 27-04-2025

```
1. salty_snack=false 14683 => pickled_vegetables=false 14552 <conf:(0.99)> lift:(1) lev:(0) [0] conv:(1)
2. salty_snack=false semi-finished_bread=false 14546 => pickled_vegetables=false 14415 <conf:(0.99)> lift:(1) lev:(-0) [0] conv:(0.99)
3. semi-finished_bread=false 14822 => pickled_vegetables=false 14688 <conf:(0.99)> lift:(1) lev:(-0) [-1] conv:(0.98)
4. salty_snack=false 14683 => semi-finished_bread=false 14546 <conf:(0.99)> lift:(1) lev:(0) [2] conv:(1.01)
5. salty_snack=false pickled_vegetables=false 14552 => semi-finished_bread=false 14415 <conf:(0.99)> lift:(1) lev:(0) [1] conv:(1)
6. pickled_vegetables=false 14830 => semi-finished_bread=false 14688 <conf:(0.99)> lift:(1) lev:(-0) [-1] conv:(0.98)
7. salty_snack=false 14683 => semi-finished_bread=false pickled_vegetables=false 14415 <conf:(0.99)> lift:(1) lev:(0) [2] conv:(1.01)
8. semi-finished_bread=false pickled_vegetables=false 14688 => salty_snack=false 14415 <conf:(0.99)> lift:(1) lev:(0) [2] conv:(1.01)
9. semi-finished_bread=false 14822 => salty_snack=false 14546 <conf:(0.98)> lift:(1) lev:(0) [2] conv:(1)
10. pickled_vegetables=false 14830 => salty_snack=false 14552 <conf:(0.98)> lift:(1) lev:(0) [0] conv:(1)
```

Applications of Apriori Algorithm:-

- Market Basket Analysis** – Finds frequently bought item pairs (e.g., bread & butter).
- Recommendation Systems** – Suggests related products (e.g., Amazon's "Frequently bought together").
- Healthcare** – Identifies disease-symptom or drug-treatment patterns.
- Fraud Detection** – Detects unusual transaction patterns.
- Web Usage Mining** – Analyzes user clickstreams for better website design.
- Inventory Management** – Optimizes stock based on purchase trends.

Used in retail, e-commerce, healthcare, banking, and more for **pattern discovery**.

Assignment-10

1)FP-Growth algorithm for association rule mining using Python libraries.

Dataset Overview

- Format: CSV (Comma-Separated Values)
- Rows: 1,4965 (based on the visible data)
- Columns: 11 (labeled 0 through 10)
- Data Type: Categorical (item names as strings)
- Purpose: Likely used for market basket analysis, such as identifying purchasing patterns, frequent itemsets, or association rules (e.g., "if a customer buys X, they are likely to buy Y").

Python code for Fp-Growth algorithm for association rule mining.

```
import pandas as pd
from mlxtend.frequent_patterns import fpgrowth, association_rules

# Step 1: Load the dataset
df = pd.read_csv('basket.csv', header=None)

# Step 2: Preprocess the dataset
# Convert the dataset into a binary (one-hot encoded) format
# Get all unique items
items = df.values.flatten()
unique_items = pd.Series(items).dropna().unique()

# Create a binary matrix (1 if item is present, 0 otherwise)
binary_data = pd.DataFrame(0, index=df.index, columns=unique_items)
for idx, row in df.iterrows():
    for item in row.dropna():
        binary_data.loc[idx, item] = 1

# Step 3: Apply FP-Growth algorithm to find frequent itemsets
# Minimum support threshold (e.g., 0.05 means itemset appears in 5% of transactions)
min_support = 0.01
frequent_itemsets = fpgrowth(binary_data, min_support=min_support, use_colnames=True)

# Step 4: Generate association rules
# Minimum confidence threshold (e.g., 0.5 means 50% confidence)
min_confidence = 0.75
rules = association_rules(frequent_itemsets, metric="confidence", min_threshold=min_confidence)

# Step 5: Sort rules by lift and display top 10
```

```

rules = rules.sort_values(by='lift', ascending=False)
print("Frequent Itemsets:")
print(frequent_itemsets)
print("\nAssociation Rules (Top 10 by Lift):")
print(rules[['antecedents', 'consequents', 'support', 'confidence', 'lift']].head(10))

# Step 6: Save results to CSV (optional)
frequent_itemsets.to_csv('frequent_itemsets.csv', index=False)
rules.to_csv('association_rules.csv', index=False)
print("\nResults saved to 'frequent_itemsets.csv' and 'association_rules.csv'.")

```

Result:-

We do not get association rules for 75 % support and 85% confidence in this data set.

->Result on lowering the support and confidence we get.

```

C:\Users\labhg\AppData\Roaming\Python\Python313\site-packages\mlxtend\frequent_patterns\fpcommon.py:161: DeprecationWarning: DataFrames with non-bool
types result in worse computational performance and their support might be discontinued in the future. Please use a DataFrame with bool type
  warnings.warn(
Frequent Itemsets:
    support           itemsets
0   0.157912      (whole milk)
1   0.051724        (pastry)
2   0.018778      (salty snack)
3   0.085873       (yogurt)
4   0.060345       (sausage)
.. ...
64  0.011160  (yogurt, whole milk)
65  0.0111628     (whole milk, soda)
66  0.013967  (rolls/buns, whole milk)
67  0.010559  (rolls/buns, other vegetables)
68  0.014836     (whole milk, other vegetables)

[69 rows x 2 columns]

Association Rules (Top 10 by Lift):
2   0.018778      (salty snack)
3   0.085873       (yogurt)
4   0.060345       (sausage)
2   0.018778      (salty snack)
2   0.018778      (salty snack)
3   0.085873       (yogurt)
4   0.060345       (sausage)
.. ...
64  0.011160  (yogurt, whole milk)
65  0.0111628     (whole milk, soda)
66  0.013967  (rolls/buns, whole milk)
67  0.010559  (rolls/buns, other vegetables)
68  0.014836     (whole milk, other vegetables)

[69 rows x 2 columns]

```

```

association_rules.csv
1 antecedents,consequents,antecedent support,consequent support,support,confidence,lift,representativity,leverage,conviction,zhangs_metric,jaccard,certainty,kulczynski
2 frozenset({'yogurt'}),frozenset({'whole milk'}),0.08587276129377172,0.15791232290831328,0.011160117615610799,0.1299610894941634,0.8229952362211855,1.0,-0.0024002495948397
3 frozenset({'rolls/buns'}),frozenset({'whole milk'}),0.10999732691793639,0.15791232290831328,0.01396685378241112,0.1269744835965978,0.8048821720437958,1.0,-0.0034030796249
4 frozenset({'other vegetables'}),frozenset({'whole milk'}),0.12209302325581395,0.15791232290831328,0.01483560545308741,0.1215106732348117,0.7694818934768153,1.0,-0.004444
5 frozenset({'soda'}),frozenset({'whole milk'}),0.097099705960973,0.15791232290831328,0.011627906976744186,0.11975223675154852,0.7583463693398951,1.0,-0.0037053331452672547
6

```

2) Fp growth algorithm for association rule mining using weka.

The basket.csv dataset is in a transactional format (items listed in columns per transaction), but Weka's Apriori algorithm requires data in a specific format, typically **ARFF (Attribute-Relation File Format)** with binary attributes (1 for item presence, 0 for absence) or a sparse format.so we need to preprocess it.

```

Python code for converting basket.csv to basket.arff.
import pandas as pd
from collections import defaultdict

# Load the dataset
df = pd.read_csv('basket.csv', header=None)

# Get all unique items
items = df.values.flatten()
unique_items = pd.Series(items).dropna().unique()

# Create a binary matrix (1 if item is present, 0 otherwise)
binary_data = pd.DataFrame(0, index=df.index, columns=unique_items)
for idx, row in df.iterrows():
    for item in row.dropna():
        binary_data.loc[idx, item] = 1

# Save to CSV (intermediate step)
binary_data.to_csv('basket_binary.csv', index=False)

# Convert to ARFF
arff_content = f"""\@RELATION market_basket

"""

for item in unique_items:
    arff_content += f"\@ATTRIBUTE {item.replace(' ', '_')} {{true, false}}\n"

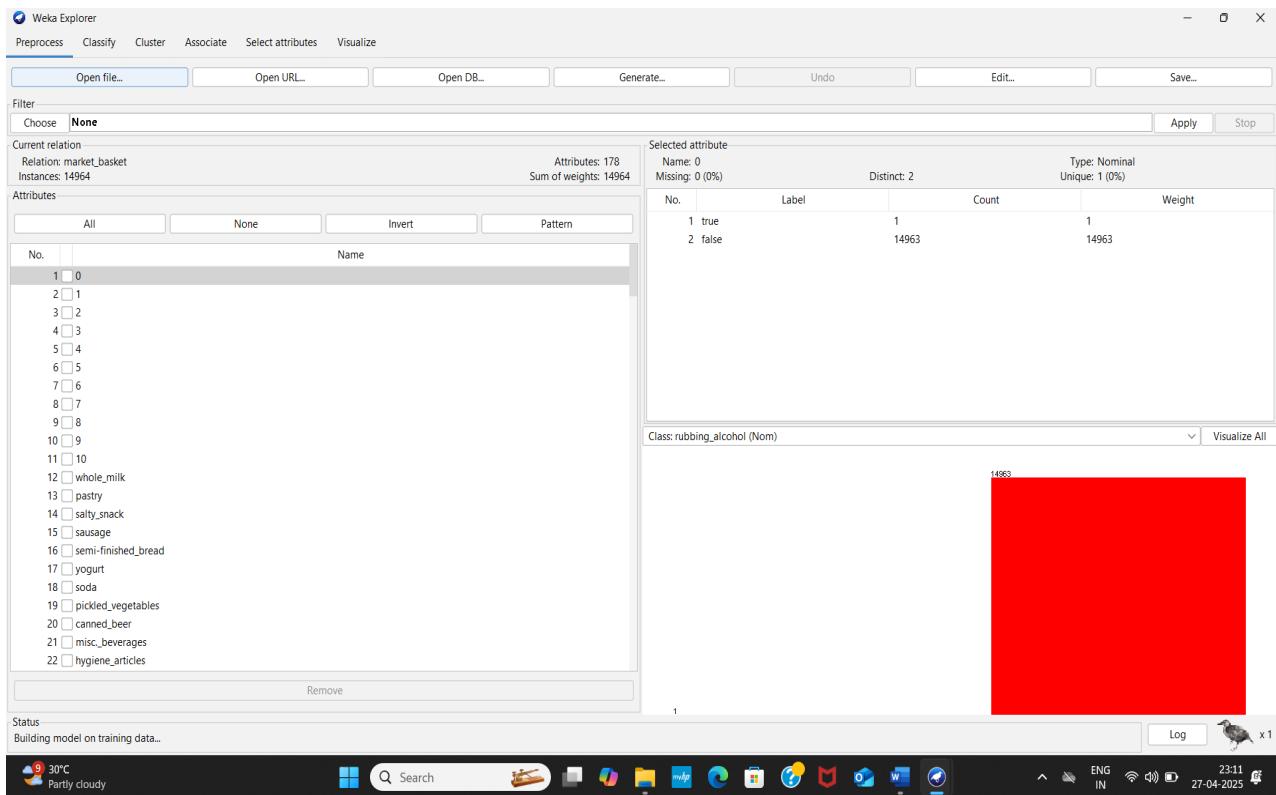
arff_content += "\n@\nDATA\n"
for _, row in binary_data.iterrows():
    row_values = ['true' if val == 1 else 'false' for val in row]
    arff_content += ','.join(row_values) + '\n'

# Write to ARFF file
with open('basket.arff', 'w') as f:
    f.write(arff_content)

print("ARFF file 'basket.arff' created.")

```

Step-2 Loading the dataset into weka



Step 2: Run the FP-Growth Algorithm

1. **Switch to Associate Tab:**
 - o Click the “Associate” tab in Weka Explorer.
2. **Select FP-Growth:**
 - o Click “Choose” and select weka.associations.FPGrowth.
 - o If FP-Growth is disabled, ensure all attributes are binary and no class attribute is set. Apply NumericToBinary or NominalToBinary filters if needed.
3. **Configure Parameters:**
 - o Click the “FPGrowth” text to open the settings window.
 - o Key parameters:
 - **lowerBoundMinSupport:** Minimum support =10%
 - **minMetric:** Minimum confidence for rules (e.g., 0.5 for 50% confidence)
 - **numRules:** Maximum number of rules (e.g., 10–100)

Step-3 Run the algorithm.

Click “Start” to execute Fp growth.

Result:-

The screenshot shows the Weka Explorer interface with the 'Associate' tab selected. The 'Associate' section is set to 'FPGrowth -P 2 -I -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1'. The 'Result list' pane displays the following information:

```

Result list (right-click for ...)
2341:11 - FPGrowth
    Start Stop
    Associate output
    === Run information ===
    Scheme: weka.associations.FPGrowth -P 2 -I -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1
    Relation: market_basket-weka.filters.unsupervised.attribute.Remove-R1-12,16,22-178
    Instances: 14964
    Attributes: 8
    pastry
    salty_snack
    sausage
    yogurt
    soda
    pickled_vegetables
    canned_beer
    misc_beverages
    === Associate model (full training set) ===
    FPGrowth found 12 rules (displaying top 10)

    1. [salty_snack=false]: 14683 ==> [pickled_vegetables=false]: 14552 <conf:(0.59)> lift:(1) lev:(0) conv:(1)
    2. [misc_beverages=false, salty_snack=false]: 14449 ==> [pickled_vegetables=false]: 14320 <conf:(0.99)> lift:(1) lev:(0) conv:(1)
    3. [misc_beverages=false]: 14728 ==> [pickled_vegetables=false]: 14596 <conf:(0.59)> lift:(1) lev:(-0) conv:(0.59)
    4. [pickled_vegetables=false]: 14830 ==> [misc_beverages=false]: 14596 <conf:(0.58)> lift:(1) lev:(-0) conv:(1)
    5. [salty_snack=false]: 14683 ==> [misc_beverages=false]: 14449 <conf:(0.98)> lift:(1) lev:(-0) conv:(0.99)
    6. [pickled_vegetables=false, salty_snack=false]: 14552 ==> [misc_beverages=false]: 14320 <conf:(0.98)> lift:(1) lev:(-0) conv:(0.98)
    7. [pickled_vegetables=false]: 14830 ==> [salty_snack=false]: 14552 <conf:(0.98)> lift:(1) lev:(0) conv:(1)
    8. [pickled_vegetables=false, misc_beverages=false]: 14596 ==> [salty_snack=false]: 14320 <conf:(0.98)> lift:(1) lev:(-0) conv:(0.99)
    9. [misc_beverages=false]: 14728 ==> [salty_snack=false]: 14449 <conf:(0.98)> lift:(1) lev:(-0) conv:(0.99)
    10. [salty_snack=false]: 14683 ==> [pickled_vegetables=false, misc_beverages=false]: 14320 <conf:(0.98)> lift:(1) lev:(-0) conv:(0.99)

```

The status bar at the bottom shows: Status OK, Log, 29°C Partly cloudy, ENG IN, 23:41, 27-04-2025.

Applications of Fp growth Algorithm:-

- Market Basket Analysis:** Finds item purchase patterns (e.g., whole milk => yogurt) for product placement, promotions, and cross-selling in retail/e-commerce.
- Recommendation Systems:** Suggests products based on user behavior (e.g., recommending rolls/buns with sausage).
- Customer Segmentation:** Groups customers by purchasing patterns for targeted marketing.
- Inventory Management:** Optimizes stock for frequently co-purchased items.
- Fraud Detection:** Identifies unusual transaction patterns in banking or insurance.
- Web Usage Mining:** Analyzes clickstream data for website optimization.
- Bioinformatics:** Discovers frequent patterns in genomic or protein sequences.
- Text Mining:** Finds co-occurring words or phrases for topic modeling or sentiment analysis.