



**Academic Year: 2022-2023**

Name:	Prerna Sunil Jadhav
Sap Id:	60004220127
Class:	T. Y. B.Tech (Computer Engineering)
Course:	Data Mining and Warehouse Laboratory
Course Code:	DJ19CEL501
Experiment No.:	01

**AIM:** Perform data Pre-processing task using Weka data mining tool

**THEORY:**

WEKA - an open source software provides tools for data preprocessing, implementation of several Machine Learning algorithms, and visualization tools so that you can develop machine learning techniques and apply them to real-world data mining problems

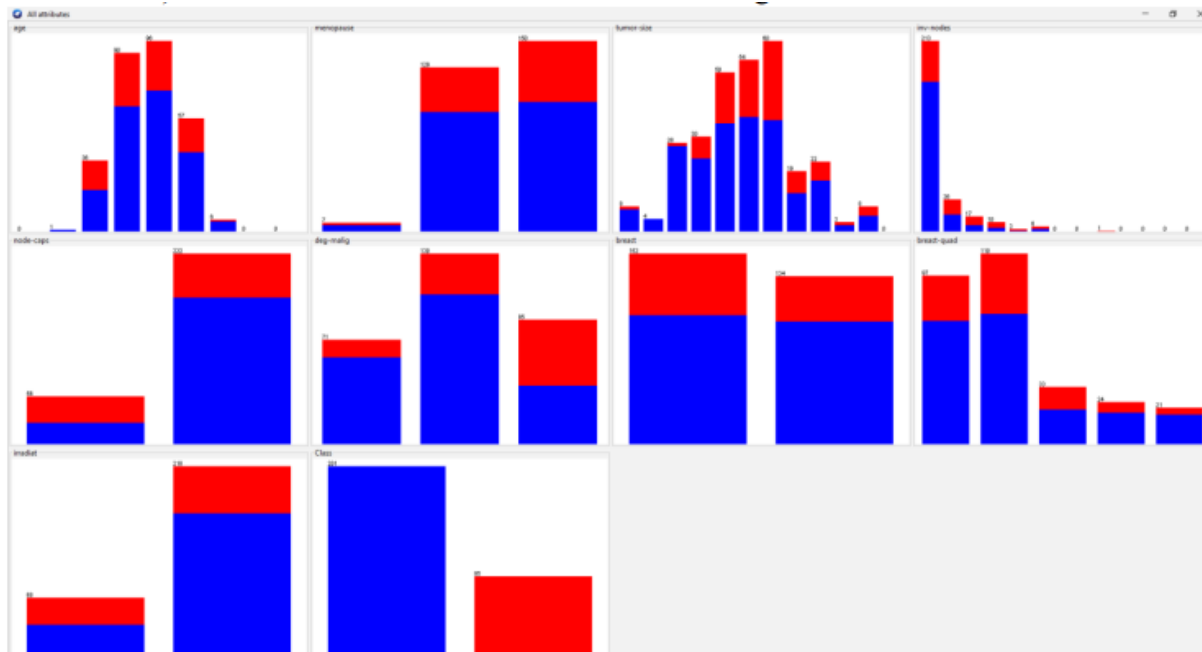
### TASKS PERFORMED THROUGH WEKA:



PREPROCESSING:

Procedure:

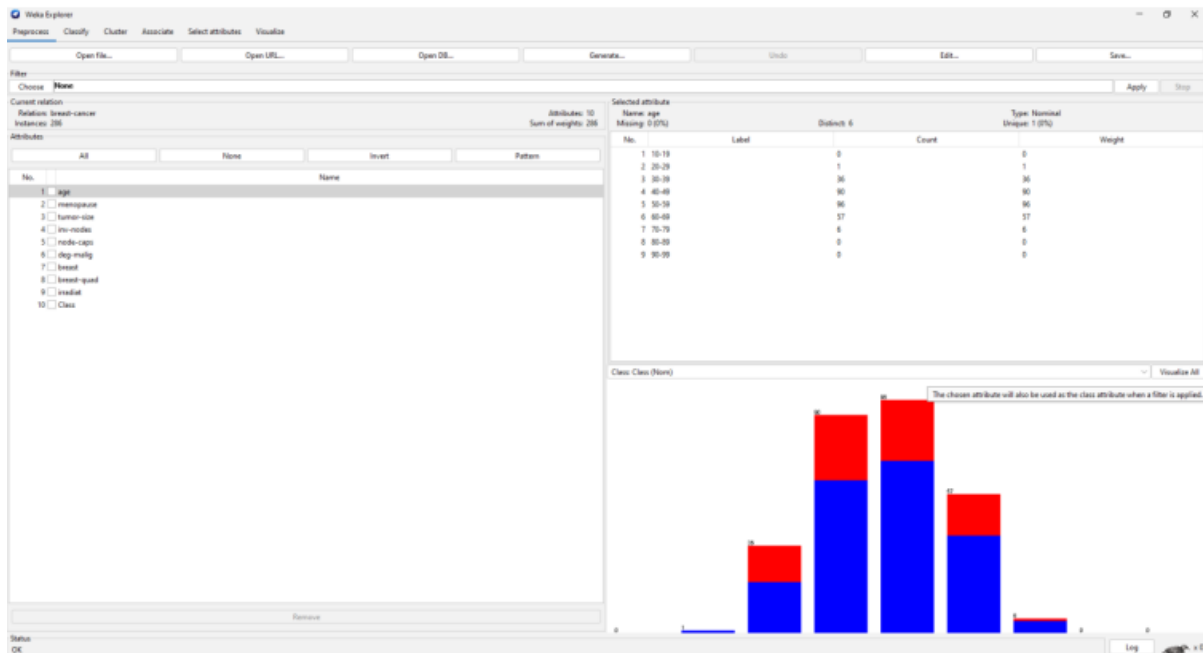
- a. **Visualize All:** Select this button to visualize histograms of all attributes.



- Filter: Choose Discretization under Unsupervised and Supervised methods. Observe the discretization and the outliers.
- IQR: Observe the IQR values for a selected attribute. Observe the outlier and extreme values
- Remove the value: Remove instances with outlier values and show the screenshots of dataset before and after the removal.

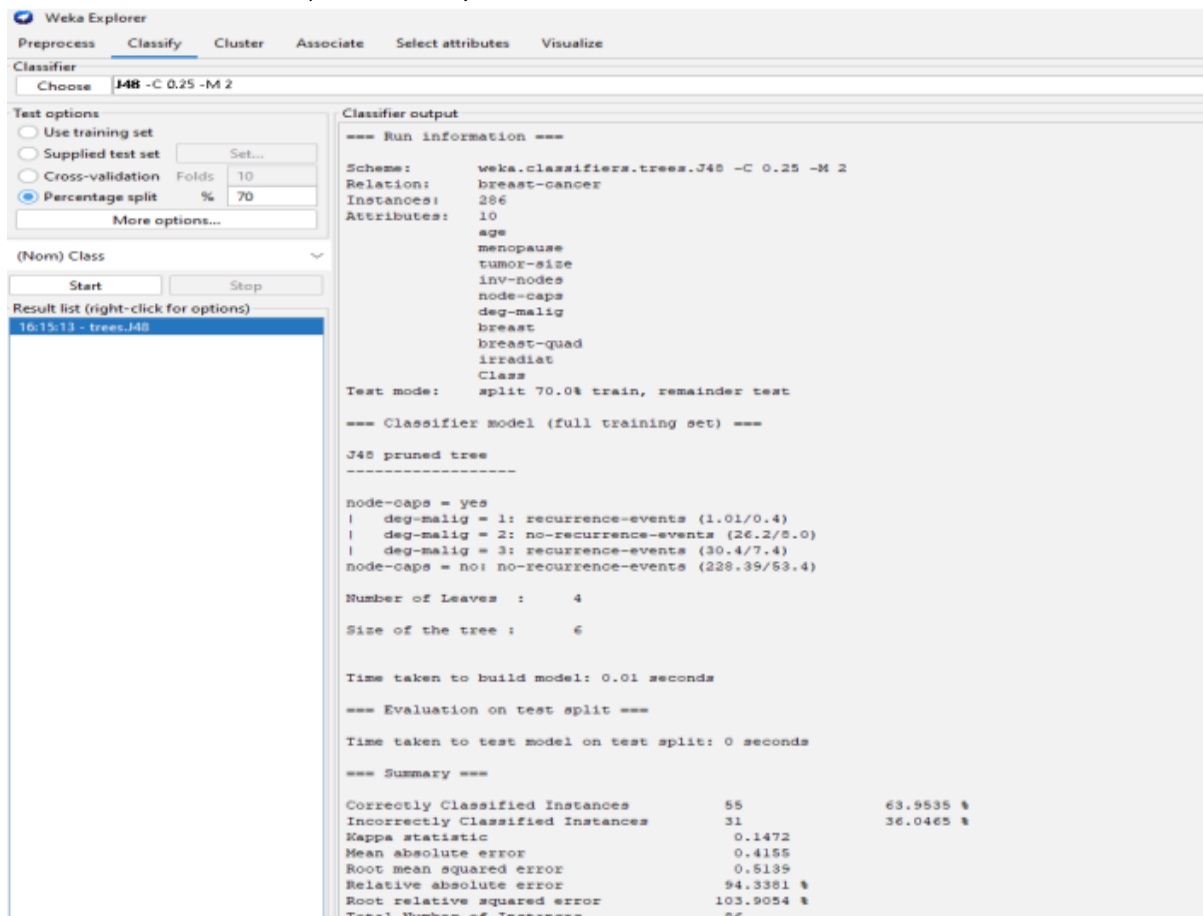


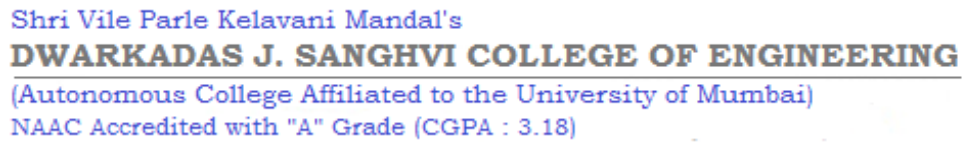
Academic Year: 2022-2023



### CLASSIFICATION:

PROCEDURE: Perform NB, kNN and DT/rule based classification





Weka Classifier Tree Visualizer: 16:18:48 - trees.J48 (breast-cancer)  
Tree View



PROCEDURE: Perform kmeans, hierarchical clustering and explain the output

The screenshot shows the Weka Explorer interface with the HierarchicalClusterer selected. The 'Cluster mode' section has 'Use training set' selected. The 'Clusterer output' section shows the command: 'Scheme: weka.clusterers.HierarchicalClusterer -B 0 -L SINGLE -P -A "weka.core.EuclideanDistance -R first-last"'. The 'Result list' shows 'Cluster 0' selected. The 'Clustered Instances' table shows 6 instances grouped into 4 clusters.

Cluster	Instances
0	279 ( 98%)
1	3 ( 1%)
2	1 ( 0%)
3	2 ( 1%)
4	1 ( 0%)



Academic Year: 2022-2023

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Cluster

Cluster mode

Use training set

Percentage split

Store clusters for visualization

Cluster output

Run information

Schema: weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -c1 -1.25 -c2 -1.0 -B 10 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10

Instances: 286

Attributes: 16

age

menopause

stage-ecog

size-codea

size-capa

deg-melig

breast-quad

irradiat

Class

Text mode: evaluate on training data

==== Clustering model (full training set) ===

KMeans

====

Number of iterations: 4

Within cluster sum of squared errors: 766.0

Initial starting points (random):

Cluster 0: 50-55,premeno,15-14,0-2,no,2,right,left\_up,no,no-recurrence-events

Cluster 1: 40-45,premeno,15-15,0-2,yes,1,right,left\_up,no,no-recurrence-events

Cluster 2: 50-55,premeno,25-25,0-2,no,1,left,left\_low,no,no-recurrence-events

Cluster 3: 50-55,pacl,40-44,0-0,yes,3,left,left\_low,yes,recurrence-events

Cluster 4: 50-55,pacl,50-54,0-2,no,1,right,right\_up,no,no-recurrence-events

Cluster 5: 50-55,pacl,20-24,0-2,no,2,left,left\_up,no,no-recurrence-events

Cluster 6: 40-45,premeno,30-34,12-14,yes,3,left,left\_up,yes,recurrence-events

Cluster 7: 40-45,pacl,30-34,0-0,no,3,left,left\_low,no,no-recurrence-events

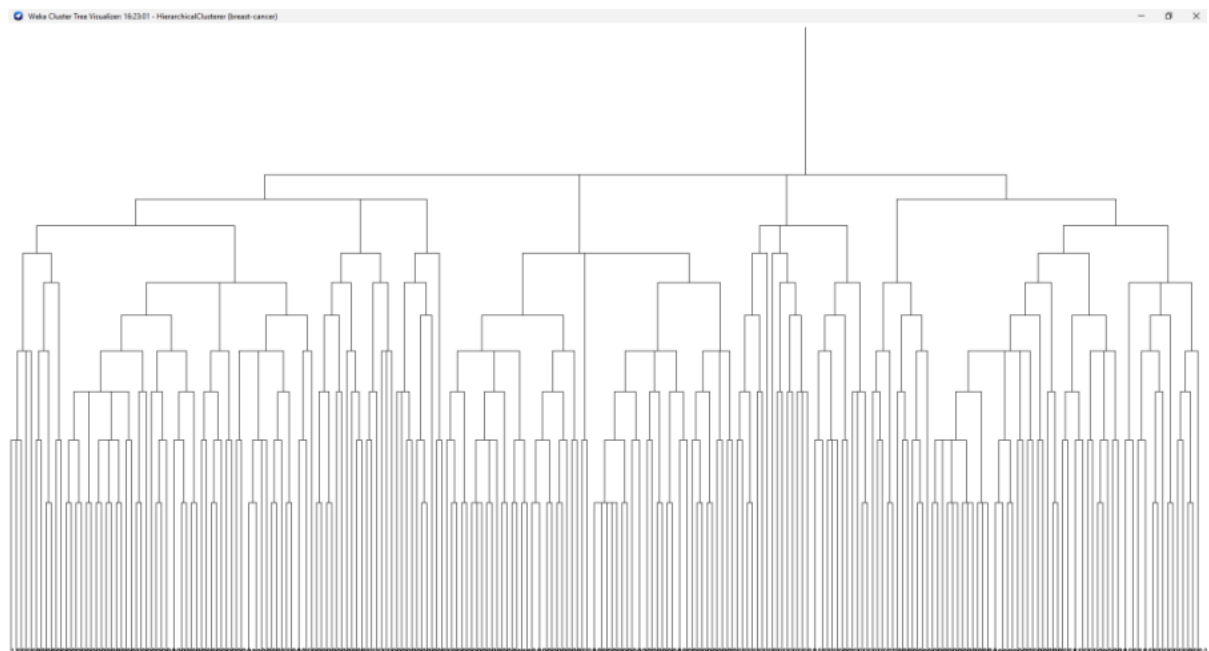
Cluster 8: 30-35,premeno,15-14,0-2,no,1,right,left\_low,no,no-recurrence-events

Cluster 9: 30-35,premeno,30-34,0-2,no,3,left,left\_up,yes,recurrence-events

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute	Full Data (286.0)	Cluster 0 (14.0)	Cluster 1 (19.0)	Cluster 2 (41.0)	Cluster 3 (20.0)	Cluster 4 (34.0)	Cluster 5 (24.0)	Cluster 6 (14.0)	Cluster 7 (29.0)
age	50-55	40-45	40-45	50-55	50-55	50-55	50-55	40-45	40-45
menopause	premeno	premeno	premeno	premeno	pacl	pacl	pacl	premeno	pacl





Academic Year: 2022-2023



### ASSOCIATION RULE:

PROCEDURE: Perform apriori algorithm and show the rules created

The screenshot shows the Weka Explorer interface with the 'Associate' tab selected. The 'Apriori' algorithm is chosen, and the 'Result list' shows two entries: '16/26/59 - Apriori' and '16/28/03 - Apriori'. The 'Apriori' output is displayed, showing the following information:

```

Scheme: weka.associations.Apriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1
Relation: supermarket
Instances: 4627
Attributes: 217
[... (list of attributes omitted) ...]
=====
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 fruit=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 945 ==> bread and cake=t 877 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(2.92)
  
```



### SELECT ATTRIBUTES:

PROCEDURE:

- Apply suitable feature selection filter like GainRatio etc. to choose relevant attributes from the list of attributes.
- Observe the ranks / priority provided by the filter.

The screenshot shows the Weka Explorer interface with the 'Select attributes' tab selected. The 'Ranker' filter is chosen, and the 'Result list' shows two entries: '16/26/59 - Ranker' and '16/28/03 - Ranker'. The 'Attribute selection output' is displayed, showing the following information:

```

Scheme: weka.attributeSelection.Ranker -T -1 -N 10000 -D 0.001 -S -1.0 -c -1
Relation: supermarket
Instances: 4627
Attributes: 217
[... (list of attributes omitted) ...]
=====
Attribute Selection output
=====
Ranker
=====
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 fruit=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 945 ==> bread and cake=t 877 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(2.92)
  
```



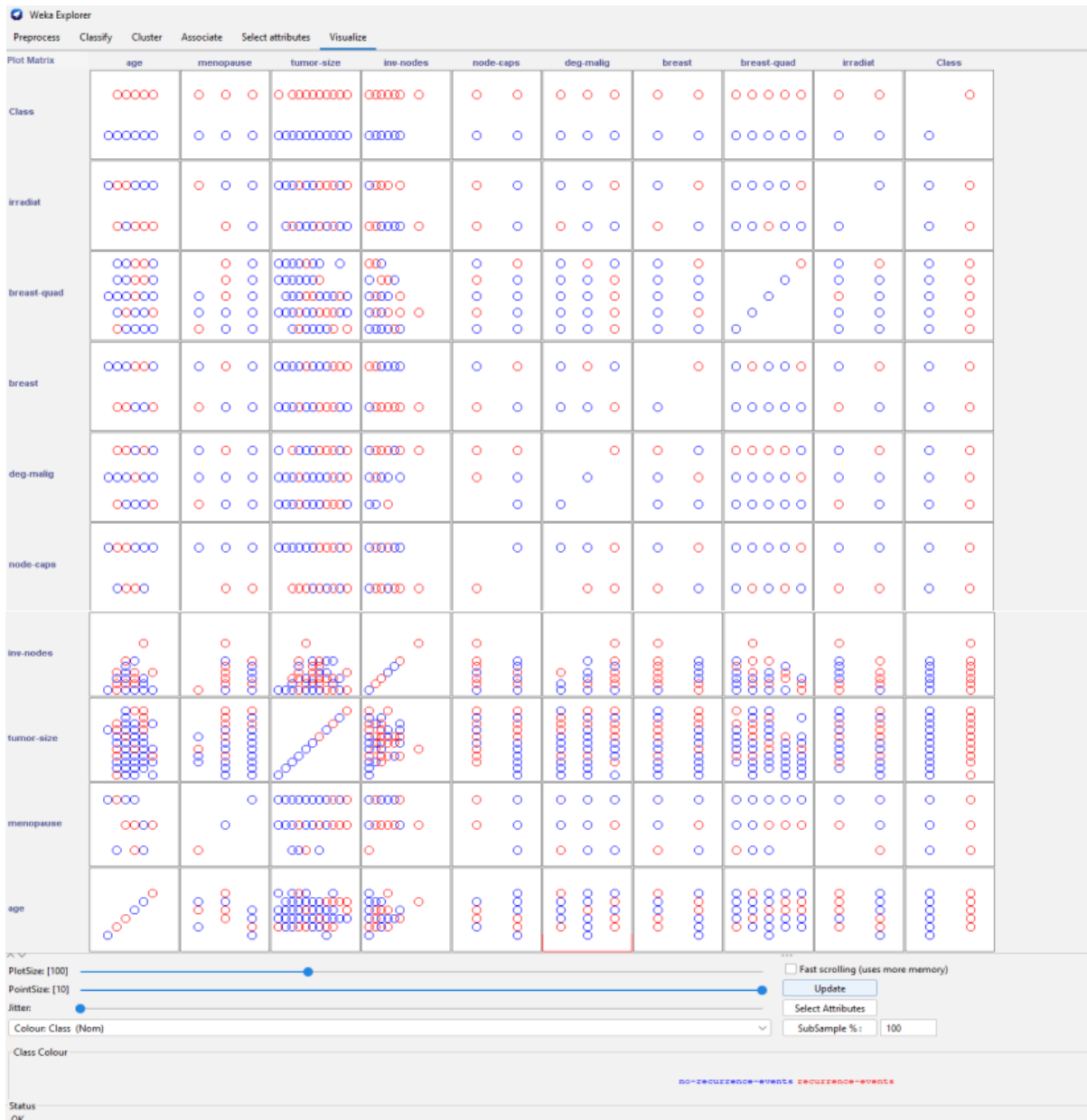
Academic Year: 2022-2023



## VISUALIZATION:

### PROCEDURE:

- Visualize scatter plot for all the attributes from a dataset selected from Weka.
- Determine correlation if any using these plots for different datasets



### CONCLUSION:



the experiment involving data pre-processing using the Weka data mining tool has been a valuable and essential step in preparing data for subsequent analysis.



Weka provides a wide range of functionalities that aid in cleaning, transforming, and organizing data, making it more suitable for data mining and machine learning tasks.