



**S. B. JAIN INSTITUTE OF TECHNOLOGY, MANAGEMENT &
RESEARCH, NAGPUR.**

Practical No. 2

Aim: Demonstration on Discretize the attribute using Weka Tool on the training data set ionosphere and compare it with different binning values and equal frequency binning.

Name of Student: Shrutika Pradeep Bagdi

Roll No.: CS22130

Semester/Year: 6th / 3rd

Academic Session: 2024-2025

Date of Performance:

Date of Submission:

AIM: Demonstration on Discretize the attribute using Weka Tool on the training data set ionosphere and compare it with different binning values and equal frequency binning.

OBJECTIVE/EXPECTED LEARNING OUTCOME:

The objectives and expected learning outcome of this practical are:

- The goal of discretization is to reduce the number of values a continuous variable assumes by grouping them into a number, b, of intervals or bins.
- In statistics and machine learning, discretization refers to the process of converting or partitioning continuous attributes, features or variables to discretized or nominal attributes /features/ variables/ intervals.

HARDWARE AND SOFTWARE REQUIREMENTS:

Hardware Requirement:

Software Requirement:

THEORY:

Discretization in Data Mining:-

Data discretization refers to a method of converting a huge number of data values into smaller ones so that the evaluation and management of data become easy. In other words, data discretization is a method of converting attributes values of continuous data into a finite set of intervals with minimum data loss. There are two forms of data discretization first is supervised discretization, and the second is unsupervised discretization. Supervised discretization refers to a method in which the class data is used. Unsupervised discretization refers to a method depending upon the way which operation proceeds. It means it works on the top-down splitting strategy and bottom-up merging strategy.

Discretization is **the process through which we can transform continuous variables, models or functions into a discrete form**. We do this by creating a set of contiguous intervals (or bins) that go across the range of our desired variable/model/function. Continuous data is Measured, while Discrete data is Counted

Two key problems in association with discretization are how to select the number of intervals or bins and how to decide on their width. Discretization can be performed with or without taking class information, , into account. These are the supervised and unsupervised ways. If class labels were known in the training data, the discretization method ought to take advantage of it, especially if the subsequently used learning algorithm for model building is supervised.

A short description of two unsupervised techniques follows.

Equal-width discretization

The algorithm first finds the minimum and maximum values of every variable, X_i , and then divides this range into a number, mX_i , of user-specified, equal-width intervals.

Equal-frequency discretization

The algorithm determines the minimum and maximum values of the variable, sorts all values in ascending order, and divides the range into a user-defined number of intervals, in such a way that every interval contains the equal number of sorted values.

Techniques of data discretization:-

Binning

Binning refers to a data smoothing technique that helps to group a huge number of continuous values into smaller values. For data discretization and the development of idea hierarchy, this technique can also be used.

Histogram analysis

Histogram refers to a plot used to represent the underlying frequency distribution of a continuous data set. Histogram assists the data inspection for data distribution. For example, Outliers, skewness representation, normal distribution representation, etc.

Cluster Analysis:

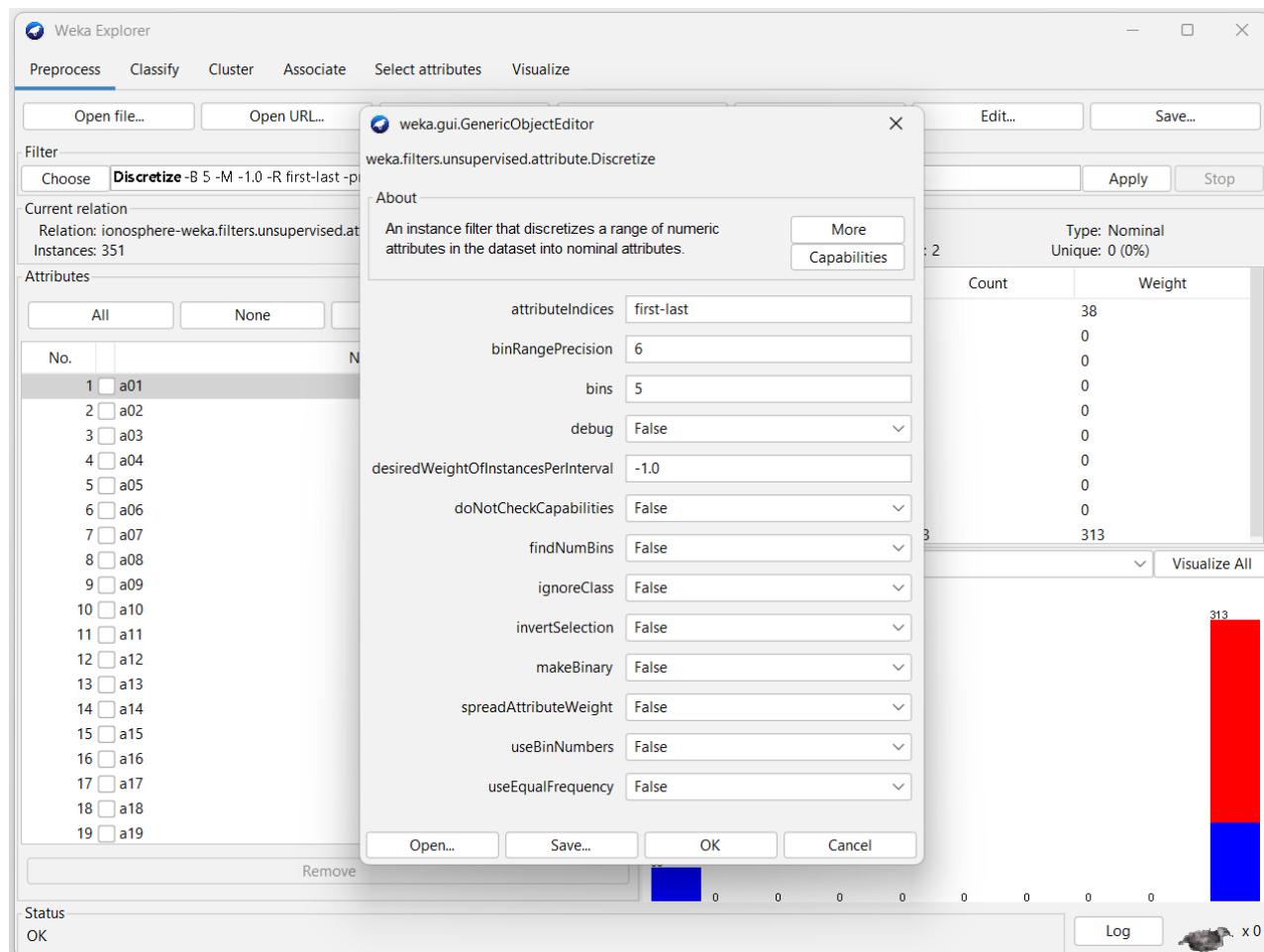
Cluster analysis is a form of data discretization. A clustering algorithm is executed by dividing the values of x numbers into clusters to isolate a computational feature of x .

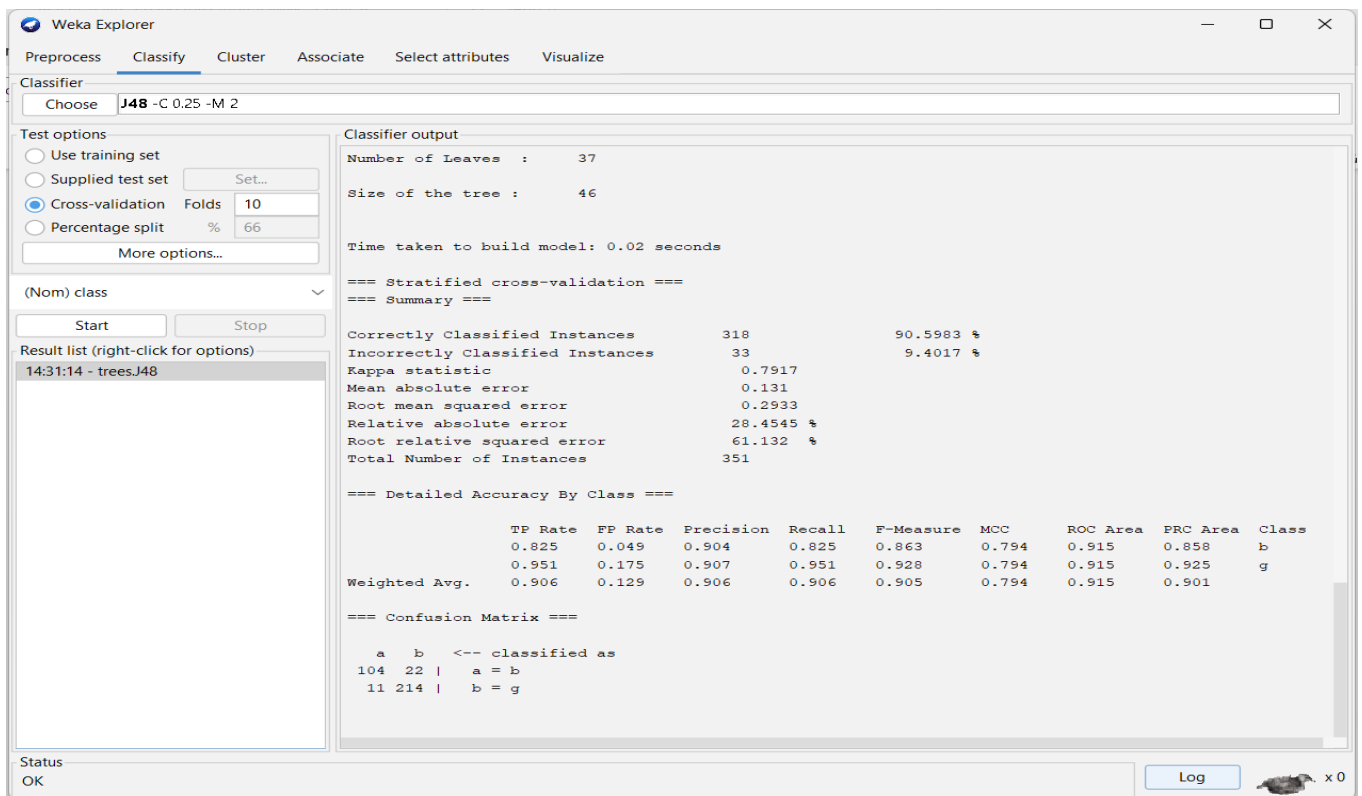
Description:

We need to use ionosphere data set which is already present in Weka training sets.

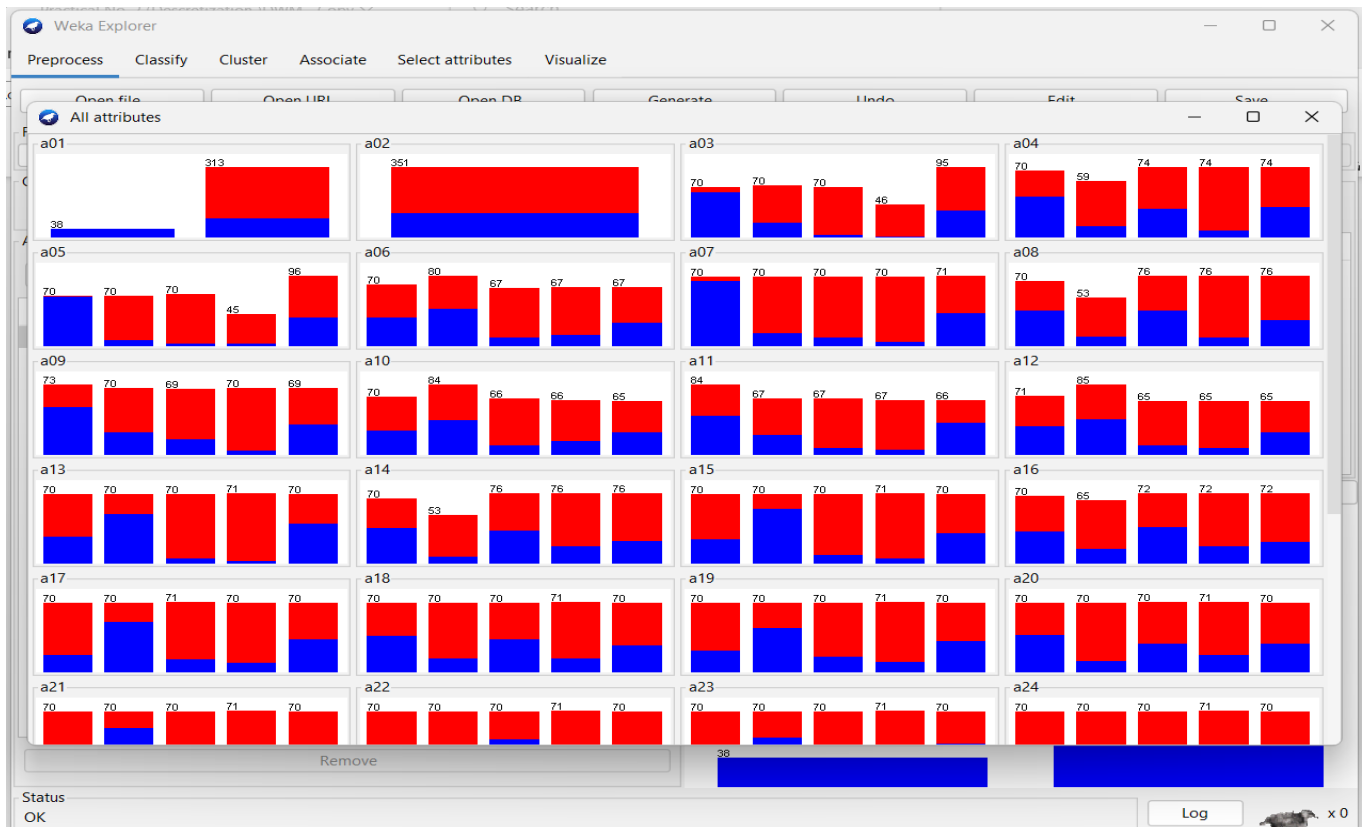
Procedure:

OUTPUT (SCREENSHOTS):





The screenshot shows the Weka Explorer interface with the 'Discretize' filter selected. The 'weka.gui.GenericObjectEditor' dialog box is open, displaying the configuration for the 'weka.filters.unsupervised.attribute.Discretize' filter. The 'attributeIndices' field is set to 'first-last', 'binRangePrecision' is 6, 'bins' is 5, 'debug' is False, 'desiredWeightOfInstancesPerInterval' is -1.0, 'doNotCheckCapabilities' is False, 'findNumBins' is False, 'ignoreClass' is False, 'invertSelection' is False, 'makeBinary' is False, 'spreadAttributeWeight' is False, 'useBinNumbers' is False, and 'useEqualFrequency' is True. The background shows the 'ionosphere' dataset with 351 instances and a list of attributes (a01 to a19). A small bar chart for attribute 'a01' is visible on the right, showing two bars: a red one with value 313 and a blue one with value 38.



The screenshot shows the Weka Explorer interface with the 'Classify' tab selected. The classifier chosen is 'J48 -C 0.25 -M 2'. The test options are set to 'Cross-validation' with 10 folds. The classifier output is displayed in the right pane, showing a stratified cross-validation summary and a detailed accuracy by class table.

Classifier output

Number of Leaves : 22
Size of the tree : 28
Time taken to build model: 0 seconds

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

Correctly Classified Instances	318	90.5983 %
Incorrectly Classified Instances	33	9.4017 %
Kappa statistic	0.7879	
Mean absolute error	0.1399	
Root mean squared error	0.2879	
Relative absolute error	30.3862 %	
Root relative squared error	60.018 %	
Total Number of Instances	351	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.786	0.027	0.943	0.786	0.857	0.795	0.893	0.858	b
	0.973	0.214	0.890	0.973	0.930	0.795	0.893	0.904	g
Weighted Avg.	0.906	0.147	0.909	0.906	0.904	0.795	0.893	0.888	

=== Confusion Matrix ===

```
a  b  <-- classified as
99  27 |  a = b
6  219 |  b = g
```

CONCLUSION:

DISCUSSION AND VIVA VOCE:

- What is Classification?
- What is Classification Accuracy?
- What is data discretization in data mining?
- Why is discretization needed in data mining?
- What is binning?
- What is the J48 Classifier?

REFERENCE:

- <http://www.cs.waikato.ac.nz/ml/weka/documentation.html>
- <http://ai.fon.bg.ac.rs/wp-content/uploads/2015/04/ML-Attribute-Discretisation-and-Selection->

Department of Computer Science & Engineering, S.B.J.I.T.M.R, Nagpur.

[Clustering-2014_eng.pdf](#)

- <https://weka.wikispaces.com/Discretizing+datasets>

Data Mining – Concepts and Techniques, Jiawei Han & Micheline Kamber, Morgan Kaufmann Publishers, Elsevier, 2nd Edition, 2006.

Observation book: (3)	Viva-Voce (3)	Quality of Submission and timely Evaluation (4)
Total:		
Sign with date:		