# An Empirical Study of Naive Bayes Classification, K-Means Clustering and Apriori Association Rule for Supermarket Dataset

Aishwarya. C. M
Department of computer science and Engineering
A.I.T College, Chikmagalur

Arpitha. D. A
Department of computer science and Engineering
A.I.T College, Chikmagalur

Prerana. B. R
Department of computer science and Engineering
A.I.T College, Chikmagalur

Rachana. C. L
Department of computer science and Engineering
A.I.T College, Chikmagalur

Karthikeyan. S. M
Asst. Prof. Department of computer science and Engineering
A.I.T College, Chikmagalur

*Abstract*- **Data mining is the process of obtaining knowledge from the large quantity of data. Generally, it is the process of interpreting data from different perspectives and summarizing it into useful information. It analyses the most useful and familiar data mining tasks like classification, clustering and association rules used by the Machine learning system, mainly in the artificial intelligent systems. The paper is mainly concerned with practical use of classification, clustering and association rule that is applied to the dataset of supermarket. In this we are studying the result of the three data mining tasks. The paper predicts the performance evaluation based on the incorrect and correct instances of data using the Naive Bayes, K-means and Apriori algorithm. The analysis of classification algorithm, clustering algorithm and association rule is done by WEKA tool.**

*Keywords- Dataset, Weka Analysis, Classification Algorithms, Clustering Algorithms, Association Rules*

## I. INTRODUCTION

Data mining is developing vastly in various fields. Data mining is adopted for several use and designed for different database [8]. The data mining task includes classification, clustering, association rule extraction, regression and visualization [7]. Classification is a data mining task that assigns items in a collection to target categories, classes. The goal of classification is to accurately obtain the target class for each case in the data [8]. Clustering is a process of partitioning a data set into set of meaningful sub-classes called clusters. Help users understand the natural grouping or structure in a data set. Association rule is procedure which is meant to find frequent patterns, correlations, associations, or casual structures from data sets found in various kinds of databases such as relational databases, transactional databases and other forms of data repositories.

For the execution of classification, clustering and association rule we have used WEKA tool. WEKA is a general collection of machine learning software written in Java, developed at the University of Waikato in New Zealand. Weka is a workbench that contains a collection of visualization tools and algorithm for data analysis and predictive modelling.

Throughout the discussion we try to understand some of the tests, analysis of classification, analysis and association rule. We make use of Weka tool to implement the different data mining tasks for the dataset supermarket. Here, we discuss about the content of data and fields which are related to the data set. The discussion is followed as first we discuss the nature of classification algorithm Naive Bayes, clustering algorithm K-means and association rule using Apriori algorithm . The next discussion is all about analysis of classification, clustering and association rule based on supermarket dataset. To analyse the supermarket datasets we use algorithms, which include Naive Bayes [4], K-means and Apriori algorithm. Finally the research results in the study of supermarket data set based on the algorithms used in the Weka tool.

## II. METHODOLOGY

We have used Weka tool for the analysis of three different data mining tasks. The algorithms of the data mining task can be directly applied to a dataset. It is also well suited for developing new machine learning schemes. The data that is used for Weka should be made into the arff format and the field should have the extension dot arff (.arff). The arff works with three sections they can be categorised into Relational section, Attribute section and the data. Whereas the data types of arff are classified into Nominal and Numeric.

### A. *WEKA GUI Chooser*

The WEKA GUI Chooser is the starting point for running the applications. There will be a choice between the command line interface (CLI), the Experiment, the Explorer and Knowledge flow. In this context we make use of explorer that gives access to all features of Weka using menu selection and form filling.

The figure1 shows the WEKA GUI chooser for commencing the main application.



Fig1: WEKA GUI chooser

### B. *Supermarket dataset*

The supermarket dataset is a dataset of point of sale information. The data is nominal and each instance represents a customer transaction at a supermarket, the product purchased and the departments involved. The performance of a comprehensive set of classification algorithms has been analysed. The dataset contains 4,627 instances and 217 attributes. The data is denormalized. Each attribute is binary and either has a value or no value. There is nominal class attribute called "total" that indicates whether transaction was less than $100 (low) or greater than $100 (high). The figure2 shows the supermarket data set opened in WEKA tool.
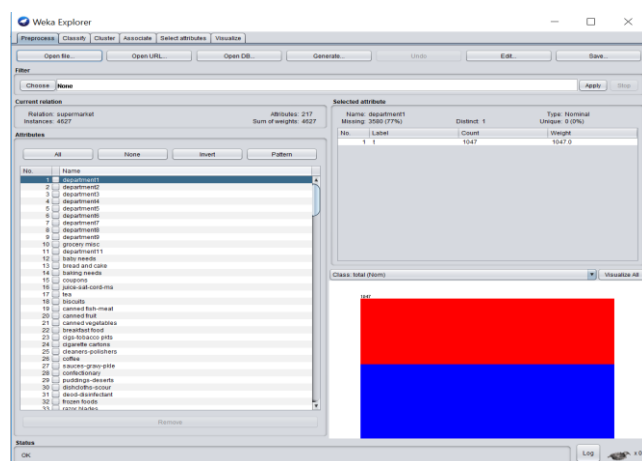


Fig2: Supermarket datasets open in WEKA

### III. DATA MINING TASK

### A. *CLASSIFICATION ALGORITHM*

There are many classification algorithms available on WEKA tool but in the paper we have selected only one classification algorithm. Analysis of classification algorithm is the assembling of data in given classes.

- *Naive Bayes*

The Naive Bayes algorithm is a simple probabilistic classifier that determines a set of possibilities by counting the constancy and combination of values in given data set [3]. This algorithm is also used in machine learning systems to conclude the new data or testing data, and it is based on the "Bayes" theory [4]. The application of this algorithm is performed by Weka tool, which provide opportunity to implement the above mentioned algorithm by using the estimator, for the numeric attributes.

### B. CLUSTERING ALGORITHM

Clustering is a technique in which a given data set is divided into groups called clusters in such a manner that the data points that are similar lie together in one cluster. Clustering plays an important role in the field of data mining due to the large amount of data set. We have used K-means clustering to obtained clustered instances.

- *K-means algorithm*

K-MEANS Clustering intends to partition n objects into k clusters in which each object belongs to the cluster with the nearest mean. This method produces exactly k different clusters of greatest possible distinction. The best number of clusters k leading to the greatest separation is not known as priori and must be computed from the data. The objective of K-Means clustering is to minimize total intra-cluster variance, or, the squared error function.

### C. ASSOCIATION RULES

This is a procedure which is meant to find the request pattern, correlations, associations or casual structures from data sets. Given a set of transactions, association rule mining aims to find the rules which enable us to predict the occurrence of a specific item based on the occurrences of the other items in the transaction.

- *Apriori algorithm*

This is the most well known association rule learning method because it may have been first and it is very efficient. In principle the algorithm is quite simple. It builds up attribute-value (item) sets that maximize the number of instances that can be explained (coverage of the dataset). The search through item space is very much similar to the problem faced with attribute selection and subset search.

### IV. SUPERMARKET DATASET ANALYSIS

- *Analysis using Naive Bayes Algorithm*

The initial analysis of this dataset is performed by default parameters that are provided by weka tool. The classification is performed by choosing "use training set".

**Special Issue - 2018**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**ICRTT - 2018 Conference Proceedings**

Which specifies that out of 4627 instances the correctly classified instances are 2948 (63.713%) and incorrectly classified instances are 1679 (36.287%). The Fig3 represents the Naïve Bayes Algorithm outcome.
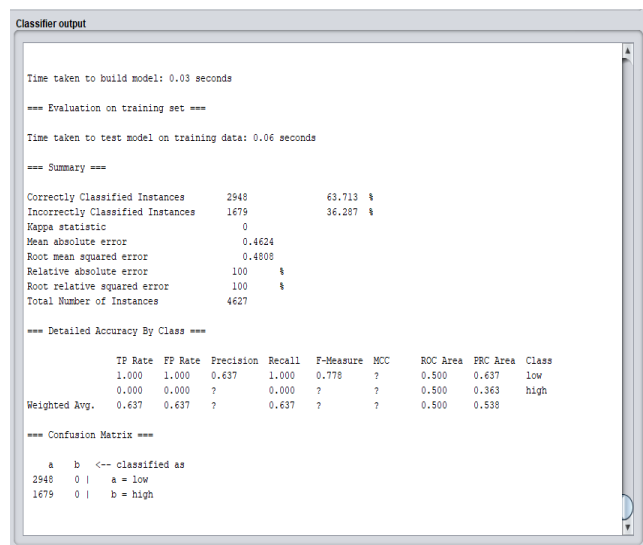


Fig3: Output of Naïve Bayes Algorithm

### B. *Analysis using K-means algorithm*

The K-means algorithm is a clustering algorithm, here training data is used for clustering it generates two clustered instances of 1679 (36%) and 2948 (64%). The total number of iterations done are two.
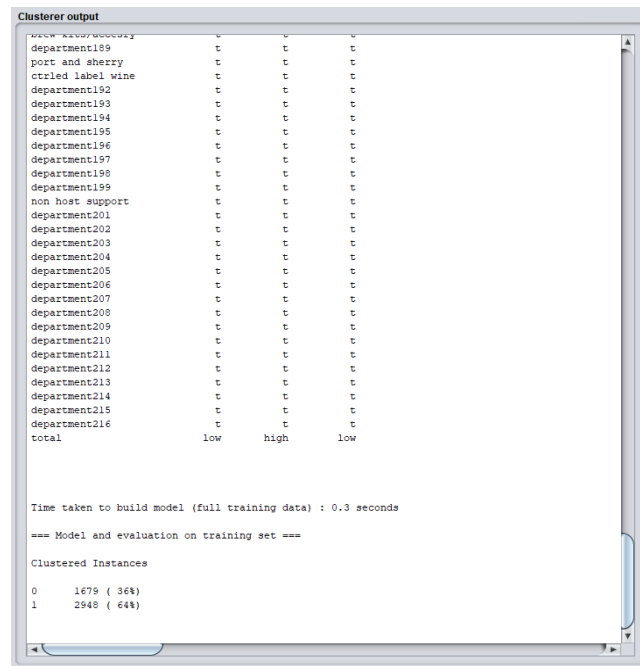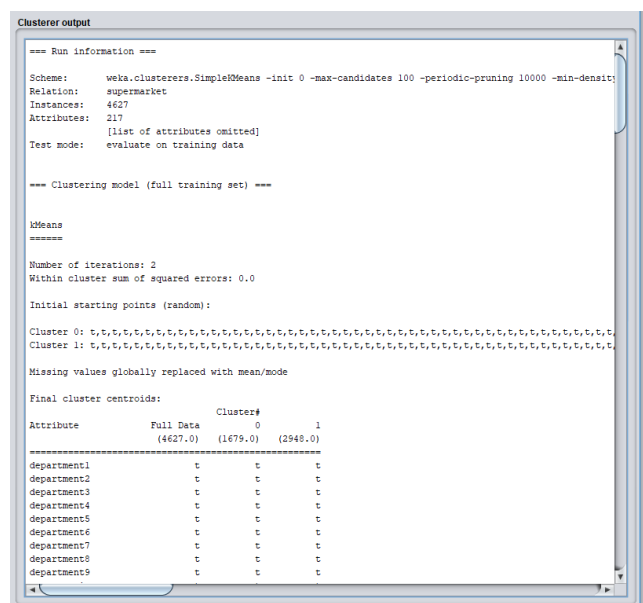




Fig4: Output of K-means Algorithm

### C. *Analysis using Apriori Algorithm*

The algorithm presented 10 rules learned from the supermarket dataset. The algorithm is configured to stop at 10 rules by default, for more rules we can configure it. Association output has no rule with coverage less than 0.91. Here are few observations:

- We can see that all presented rules have a consequent of "bread and cake".
- All presented rules indicate a high total transaction amount.
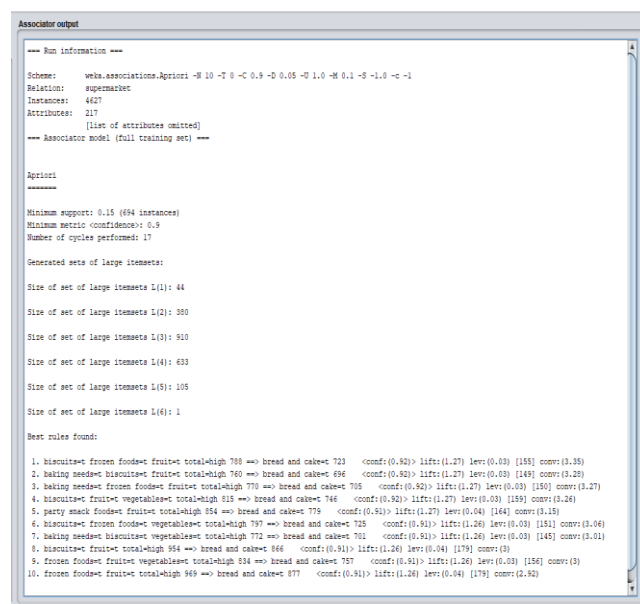- "biscuits" an "frozen foods" appear in many of the presented rules.



Fig5: Output of Apriori Algorithm

**Special Issue - 2018**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**ICRTT - 2018 Conference Proceedings**

## V.    CONCLUSION

The purpose of a data mining tasks is normally either to create a descriptive model or a predictive model. This approach analyses the application of three different algorithms of data mining task in the data set. Naive Bayes algorithm is intended primarily for the work with nominal attributes. The performance depends on the classification algorithm that is adopted. The K-means algorithm results in 2 cluster instances and the Apriori algorithm gives us the information that if itemset is frequent then all its subset is frequent and the generates several rules. It is more efficient to use an algorithm like Apriori rather than deducing rules by hand.

## ACKNOWLEDGEMENT

## REFERENCES

[1] Thangaraju, R. Mehla,  Analysis of  KStar Classifier over Liver Disease, International Journal of Advanced Research in Computer science Engineering & Technology (IJARCET) Volume 4 Issue 7,July 2015

[2] Ho Tin Kham. "Random subspace Methods for Constructing Decision Forest"

[3] George Dimytoglou, James Adam, and Carol M. Jhim, "Comparisons of C4.5 and Naive Bayes Classification for the Analysis of Lung Cancer Survivabilities"

[4] Olivier C. Fhran, kois and Philip Leray Study of the Tree Augmented Naive Bayes Classification from deficient datasets LITIS. Sain-Etienne-De-Rouary, France.

[5] Jangtao Ron*,Sau Dan Le, Xianlo Chn , Ben Ka, Renold Chenk and David Cheunk   Naive Bayer Classification of incalculable Data Department of Computer Engineering, Son Yaat-son University, Guangzhou, China

[6] Lio Bhreman, Jerom Fridman, Richerd Olshan, Charle Stone "Regression Tree" (Wardsworth).

[7] Ghopi Gandi, Rohith Shreevastav Modified k-methods algorithms for analysis and application to increase scalabilities and efficiencies for larger datasets.