# AN IN-DEPTH ANALYSIS OF THE SELECTED CLASSIFICATION TECHNIQUES OF DATA MINING

**Bhavay Bajaj**

*Gateway International School*

## ABSTRACT

*Classification is an information mining task that doles out things in an assortment to target classifications or classes. The extent of grouping is to precisely anticipate the objective type for each case in the information. In theory, construct preparing system, a grouping calculation discovers connections between the value of the indicators and the upsides of the objective. Distinctive arrangement calculations utilize different strategies for finding relationships. These connections are summed up in a model, which later is applied to an alternate informational collection where the class tasks are obscure. The arrangement has numerous applications in client division, business displaying, showcasing, credit examination, bio clinical and drug reaction demonstrating. This paper presents the review and analysis of five arrangement calculations physically Bayesian organization, j48, strategic model tree, irregular tree and rep tree for liver issues dataset and the presentation of these calculations are thought about utilizing the different exhibition measurements, for example, Precision, Recall and F measure in which arbitrary tree calculation gives 100% exactness. The exploratory outcome shows that the random tree gives high precision than the Bayesian calculation, j48, strategic model tree and rep tree.*

## I. INTRODUCTION

Information mining is the example of analysing huge prior data set to produce new information. The data acquired from information mining is ideally new and useful. , generally speaking, the objective of the information mining measure is to extricate data from an informational index and altered it into a justifiable construction for additional utilization. The crude investigation step includes data set and information the board highlights, information pre-processing, model and suspicion contemplations, intriguing quality measurements, muddled contemplations, post handling of noticed designs, representation and internet refreshing. Information mining is the examination step of the KDD(knowledge disclosure in the data set) measure [1]. The simple information mining task is the self-loader or programmed examination of huge amounts of information to separate earlier unseen, fascinating examples, for example, gatherings of information records, surprising records and conditions. Information handling depends on complex calculations that consider information division to recognize illustrations and patterns, identify peculiarities, and foresee the likelihood of different situational results.

New patterns are Distributed Data Mining, Multimedia Data Mining, Spatial and Geographic DataMining, Ubiquitous Data Mining, Time Series and Sequence Data Mining, Application Exploration, versatile and intelligent information mining strategies, reconciliation of information mining with data set frameworks, information stockroom frameworks and web data set frameworks, Visual information mining, New techniques for mining complex kinds of information, Biological information mining, Data mining and programming, Web mining, Real-time information mining and Privacy insurance and data security in information mining [2]. Information mining measure

35

applications are life sciences, client relationship the board, web applications, producing, serious advantage, insight, retail, finance, banking, PC, organization, security, observing, observation, showing support, environment demonstrating, stargazing financial information examination, retail industry, media transmission industry, natural information investigation, other logical applications and interruption revelation. Information mining benefits are Marketing, Retail, finance, banking, assembling and legislatures. Information mining drawbacks are protection issues, security issues, Misuse of data and inaccurate data. Characterization is the course of consequently making a model of classes from many records containing class marks.

Part 1 examines the presentation of information mining and the order. Segment 2 gives a brief clarify writing survey. Section 3 explains the strategy and utilized in order calculations. Segment 4 conversation results are described. Segment 5 and finishes up this examination work.
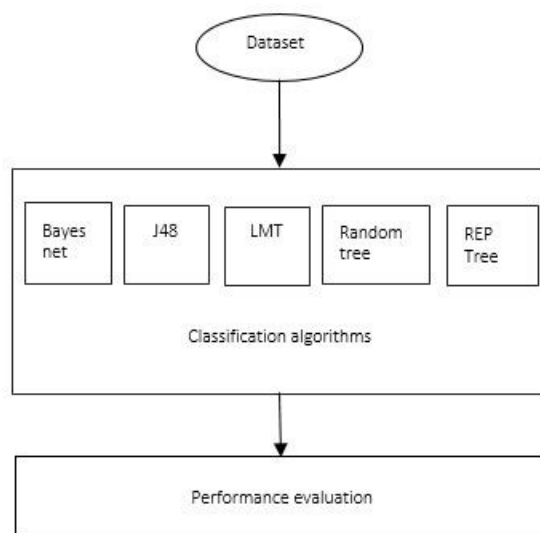
## II. PROCEDURE



Fig 1: Proposed Work Procedure

Figure 1 shows the procedure of the proposed work that comprises two stages to be specific characterization stage and the execution assessment stage. The liver problems dataset is picked as a trial dataset. The characterization stage utilizes five arrangement calculations, specifically Bayes organization, j48, logistic model tree, random tree and rep tree. These calculations are dissected and approved using the diverse exhibition assessment measurements.

### 1. Arrangement Algorithms

Grouping is the information mining task that assigns all records in the informational index to one of only a handful of exceptional predefined classes. The informative collection is partitioned into preparing and test sets. The preparing set has realized class marks while the test set names are obscure. The most famous and normal are probably adjusted and introduced here independent of their abilities, effortlessness, and strength. Order is in the like manner portrayed as the errand of target work learning

36

for planning each characteristic set to its relating, class name. There are various grouping calculations, for example, Bayes network,j48, logistic model tree, random tree and rep tree [13].

### 1.1BayesNetwork:

The bayesian organization, likewise called conviction organizations, is a graphical model for likelihood connections among a bunch of factors includes; this Bayesian organization comprises two parts. The first part is primarily a coordinated non-cyclic diagram (DAG). The hubs in the chart are known as the irregular factors, and the edges between the hubs address the probabilistic conditions among the relating arbitrary factors. The second part is a bunch of boundaries that portray the contingent likelihood of every element taken its folks. A Bayesian organization describes a framework by determining connections of contingent reliance between its variables. The restrictive conditions are addressed by a coordinated non-cyclic diagram, wherein every hub [14].

### 1.2 J48:

J48 is an open-source Java execution of the C4.5 algorithm.C4.5 is a program that makes a choice tree dependent on many named input information. It is a straightforward C4.5 choice tree for characterization; it makes a paired tree. The choice tree approach is the most helpful in the characterization issue. With this strategy, a tree is worked to display the arrangement cycle. When the tree is fabricated, it is useful to each tuple in the data set and results in arrangement; however, the tree is assembled. It overlooks the missing qualities. The fundamental thought of J48 is to partition the information into range dependent on the property estimations for what is found in the preparation Sample. It permits characterization either as a decision tree or rules made dependent on the test set gave [15].

### 1.3 Logistic Model Tree:

It is a characterization model with a related preparing calculation that combines strategic relapse (LR) and choice tree learning [16]. It is likewise called a rationale model, which is utilized to demonstrate the dichotomous result of factors. Calculated model trees depend on the previous idea of a model tree: a choice tree with direct relapse models at its leaves to give a piecewise straight relapse model.

### 1.4 Random Tree:

It is a directed classifier; it is a group learning calculation that produces numerous singular students. It utilizes a stowing thought to create an irregular arrangement of information for building a choice tree; an arbitrary tree is an assortment of tree indicators called timberland [17].

### 1.5 Rep Tree:

REP implies decreased blunder pruning. REP Tree is a rapidly choice tree student which constructs a choice and relapse tree utilizing data gain as the splitting measure and prunes it using decreased error snipping [18].

# III. RESULTS AND DISCUSSION

## A. Dataset Description

Table 1 shows the liver problems informational collection with 345 occurrences. The accompanying ascribes mcv, alkphos, sgpt, sqot, gammagt, beverages, and choice have been utilized to investigate liver issues information because of its sickness capability [19].

Table 1 Data set of Liver Disease

| ATTRIBUTE | DESCRIPTION |
|---|---|
| MCV | Mean corpuscular volume |
| ALKPHOS | Alkaline phosphotase |
| SGPT | Alamine aminotransferase |
| SQOT | Aspartate aminotransferase |
| GAMMAGT | Gamma-glutamyl transpeptidase |
| DRINKS | Number of half-print equivalents of alcoholic beverages drunk per day |
| SELECTOR | Selector field used to split data into two sets |

Table 2 shows that the five arrangement models, created with the chose information mining calculations, are looked at by utilizing the accompanying assessment measures: % of accurately and inaccurately grouped cases and kappa measurement. These are extraordinary measures for the assessment of information digging models for order. Kappa measurement is a proportion of non-arbitrary understanding among onlookers and estimation of a similar clear-cut variable. When looking at the time taken for various calculations, the Bayes net and arbitrary tree take the most un-processing time. Figure 2 shows the accuracy and mistakenly characterized cases from the aftereffect of the five order calculations and saw that random tree calculations give preferred grouping results over different estimates.

Table 2:  Classifiers Output

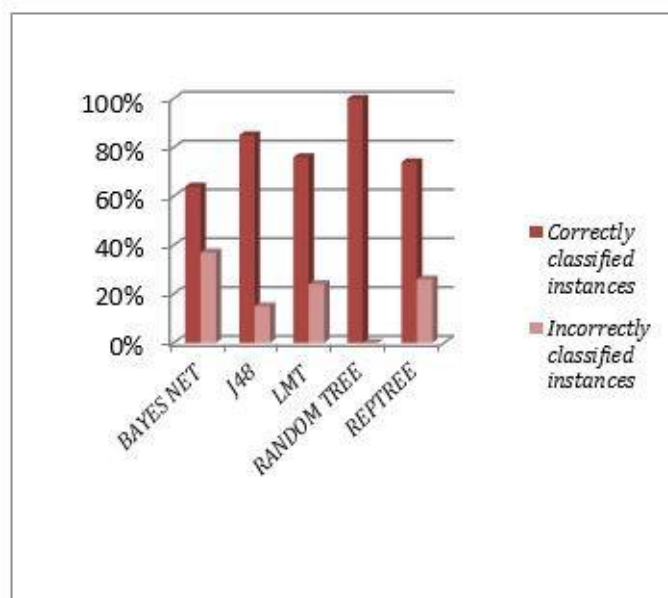| Evaluation criteria | Classification Algorithms | | | | |
|---|---|---|---|---|---|
| | Bayes net | J48 | LMT | Random tree | Rep tree |
| Correctly classified instances | 218 | 292 | 261 | 345 | 256 |
| | 64% | 85% | 76% | 100% | 74% |
| Incorrectly classified instances | 127 | 53 | 84 | 0 | 89 |
| | 36% | 15% | 24% | 0% | 26% |
| Kappa statistics | 0.24 | 0.67 | 0.48 | 1 | 0.43 |
| Time taken sec | 0.02 | 0.03 | 0.53 | 0.02 | 0.03 |
| Accuracy | 64% | 85% | 76% | 100% | 74% |



Figure 2. Correctly and Incorrectly Classified Algorithms

Table 3. Accuracy Results

| Classifier | Precision | Recall | Fmeasure |
|---|---|---|---|
| BAYES NET | 0.63 | 0.63 | 0.63 |
| J48 | 0.85 | 0.84 | 0.84 |
| LMT | 0.75 | 0.75 | 0.75 |
| RANDOM TREE | 1 | 1 | 1 |
| REO TREE | 0.76 | 0.74 | 0.72 |

The five arrangement calculations have been approved utilizing the five important measurements like accuracy, review and F-measure, and these approval results are displayed in Table 3. The same is portrayed in fig.3. The outcome shows that the random tree classifier gives preferable arrangement exactness as 1.0 over the other four calculations.
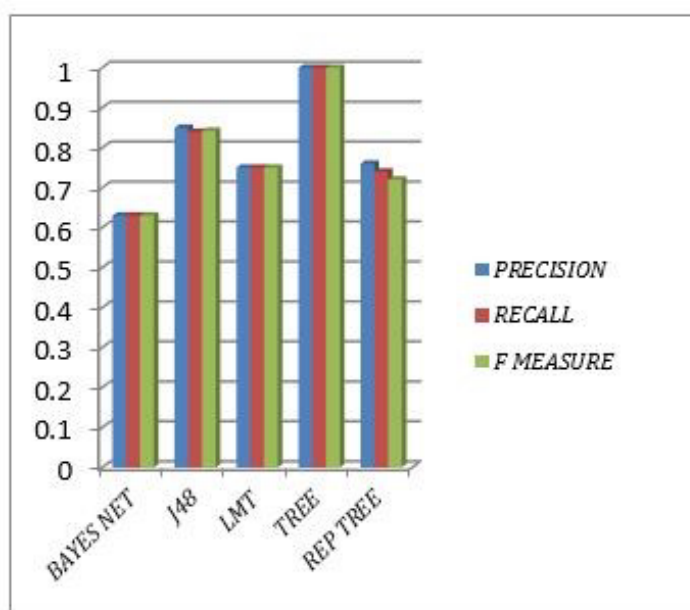


Figure 3. Performance Evaluation

## IV. CONCLUSION

This paper directed a broad review on five characterization calculations, and the trial result shows that the irregular tree classifier gives better precision of 100%, which requires 0.02 seconds for preparation. The second-best calculation is the j48, which gives exactness 85% and needs 0.03 seconds then the

40

arbitrary tree calculation. The third best calculation is a strategic model tree which provides an accuracy of 76% and requires 0.53 seconds. The fourth best calculation is the rep tree which gives exactness 74% and requires 0.03 seconds for preparing. At long last, the Bayes network offers the least precision, 64% and requires 0.02 seconds for preparing the cases.

# REFERENCES

[1]. Fayyad, Ussama; Piattetsky-Shapiro,Gregory; Smyth, Padhraic (1996). "From Data Mining to Knowledge Discovery in Databases" (PDF). Retrieve December 2008.

[2].          https://graduatedegrees.online.njit.edu/resources/mscs/mscs-articles/current-trends-in-data-mining/

[3]. Yun Wan, Dr. QigangGao," An Ensemble Sentiment Classification System of Twitter Data for Airline Services Analysis",2015 IEEE 15th International Conference on Data Mining Workshops.

[4]. Ibrahim M. El-Hasnony, Hazem M. El-Bakry, Ahmed A. Saleh, ," Classification of Breast Cancer Using Softcomputing Techniques", International Journal of Electronics and Information Engineering, Vol.4, No.1, PP.45-54, Mar. 2016.

[5]. Anita kumar,"A Study on Cancer Perpetuation Using the Classification Algorithms", International Journal of Recent Research in Mathematics Computer Science and Information Technology Vol. 2, Issue 1, pp: (96-99), Month: April 2015 – September 2015, Available at: www.paperpublications.org

[6]. HakizimanaLeopord, Dr. Wilson KiprutoCheruiyot, Dr. Stephen Kimani," A Survey and Analysis on Classification and Regression Data Mining Techniques for Diseases Outbreak Prediction in Datasets", The International Journal Of Engineering And Science (IJES) || Volume || 5 || Issue || 9 || Pages || PP -01-11 || 2016 || ISSN (e): 2319 – 1813 ISSN (p): 2319 – 1805.

[7]. NityaUpadhyay, VinodiniKatiyar," A Survey on the Classification Techniques In Educational Data Mining", International Journal of Computer Applications Technology and ResearchVolume 3– Issue 11, 725 - 728, 2014, ISSN: 2319–8656.

[8]. Fabien Lotte, Marco Congedo, Anatole Lecuyer, FabriceLamarche, Bruno Arnaldi, "A review of classification algorithms for EEG-based Brain computer interfaces", Journal of Neural Engineering, IOP Publishing, 2007, 4, pp.24. <inria-00134950>.

[9]. Patel Pinky S. Devendra V. Thakor," A Survey of Email Classification Algorithms in Data Mining", International Journal of Engineering Technology, Management and Applied Sciences www.ijetmas.com January 2015, Volume 3 Issue 1, ISSN 2349-4476.

[10]. Arvind Kumar, ParminderKaur andPratibha Sharma,"A Survey on Hoeffding Tree Stream

Data ClassificationAlgorithms", CPUH-Research Journal: 2015, 1(2), 28-32ISSN (Online): 2455-6076 http://www.cpuh.in/academics/academic_journals.php

[11]. VandanaKorde and C NamrataMahender," Text classification and classifiers:A survey", International Journal of Artificial Intelligence & Applications (IJAIA), Vol.3, No.2, March 2012.

[12]. Divya Jain, VijendraSingh ," Utilization of Data Mining Classification Approach for Disease Prediction: A Survey", I.J. Education and Management Engineering, 2016, 6, 45-52 Published Online November 2016 in MECS (http://www.mecs-press.net) DOI:10.5815/ijeme.2016.06.05.

[13]. L. Tao, F. Sun, and S. Yang, A fast and robust sparse approach for hyper spectral data classification using a few labelled samples," IEEE Transactionson Geoscience and Remote Sensing, vol. 50, no. 6, pp. 2287-2302, 2012.

[14]. Delveen Luqman Abd AL-Nabi1, Shereen Shukri Ahmed2, Survey on Classification Algorithms for Data Mining:(Comparison and Evaluation),Computer Engineering and Intelligent SystemsISSN 2222-1719 (Paper) , ISSN 2222-2863 (Online)Vol.4, No.8, 2013.

[15]. R. Sivanesan1, K. Devika Rani Dhivya2 , "A Review on Diabetes Mellitus diagnoses using classification on Pima Indian Diabetes Data Set", International Journal of Advance Research in Computer Science and Management Studies Research Article / Survey Paper / Case Study, Volume 5, Issue 1, January 2017, Available online at: www.ijarcsms.com

[16]. Wikipedia (2017). Logistic model tree. Available:https://en.wikipedia.org/wiki/Logistic model tree.

[17]. Jiawei Han and Micheline Kamber Data Mining:Concepts and Techniques, second edition.

[18]. Ghosh, S. R. and Waheed, S. (2017). Analysis of classification algorithms for liver disease.

diagnosis. Journal of Science, Technology and Environment Informatics, 05(01), 361-370. https://doi.org/10.18801/jstei.050117.38 .

[19]. C.L. Blake, D.J. Newman, S. Hettich and C.J. Merz. (2012) UCI machine learning repository databases. [Online]. Available: http://mlr.cs.umass.edu/ml/machine-learning-databases/0022

[20]. Cheng, Hong, et al. "Discriminative frequent pattern analysis for effective classification." Data Engineering, 2007. ICDE 2007. IEEE 23rd International Conference on. IEEE, 2007.

[21]. M. El-Hasnony, H. M. El Bakry, A. A. Saleh, "Comparative study among data reduction techniques over classification accuracy," International Journal of Computer Applications, vol. 122, no. 2, pp. 8,15, 2015.

[22]. John C. Bailar, Thomas A. Louis, Philip W. Lavori, Marcia Polansky, "A Classification for Biomedical Research Reports," N Engl J Med, Vol. 311, No. 23 pp. 1482-1487, in the year 2010.

[23]. Ada, Rajneet Kaur. "Using Some Data Mining Techniques to Predict the Survival Year of Lung Cancer Patient." (2013).

[24]. M. A. Nishara Banu, B. Gomathy ,"Disease Forecasting System Using Data Mining Methods", 2014 International Conference on Intelligent Computing Applications.