

CLASSIFICATION TECHNIQUES OF DATA MINING: A REVIEW

Pooja Arora¹ and Ajay Choudhary²

¹Department of Information Technology, Poddar International College, University of Rajasthan, Jaipur, 302004, India

²Department of Information Technology, Poddar Management And Technical Campus, Rajasthan Technical University, Jaipur, 302004, India

ABSTRACT

Characteristically, Data mining is a set of methods for computerized discovery of unidentified, suitable, narrative, constructive and explicable prototypes from huge databases which actively employ in practice of decision making. Habitually business intelligence association and economic forecaster adopt this process but with time span it progressively more utilize in the sciences for takeout only imperative data from huge data sets created by modern investigational and observational process. There a number of data mining-based methods are available for aid to professionals and each and every model perform differently, suitability of model depends on the base of function and their attributes types. Additionally, almost of the offered method are regulate to drive with an explicit type of statistics that outperformed with one specific application but struggle with other test data set. Therefore, design of competent model is still a current challenge of related field. This paper presents the modern status, challenges and the present hitches of accessible classification procedures of data mining.

Keyword: Data Mining, Classification Techniques, Naïve Bayes, Decision Tree

INTRODUCTION

Today, due to promising results the practice of data mining has receives quite big attention in each and every working field. Over past decades, much of investigators indicated that data mining techniques has offers great potential benefits for decision-making system. In modern era, with the advancement of techniques and digitalization of information the experts of divers working fields employing a variety of data mining procedures to manage data with an aim to improve quality of decision making. On the other hand, experts use data mining practices for advanced statistical analysis to find out relationships between features and hidden knowledge within the composed datasets which provides useful result for future assistance [1]. Classification is a practice to choosing hypothesis from an alternative set. Classically, classification practices utilize some pre-set dataset recognized as training set to make familiarity for forecasting of novel instances from assorted data [2-4].

Elementary Modules of Classification Procedures

Classification is learning processes of digital device

To forecast naïve instances the machine has to be trained with some assorted dataset, in modest form the process of learning through some assorted set of data is recognized as classification. Such process also recognized as supervised process of machine learning that has started with the assortment of statistics from several sources. During learning phase, a classifier practice built portraying set of previously determined classes that will portrayed in future as rules [5]. There are excessive learning practices is in real world, which can be estranged into a widespread category. The figure 3.4 illustrate four foremost components of classification process.

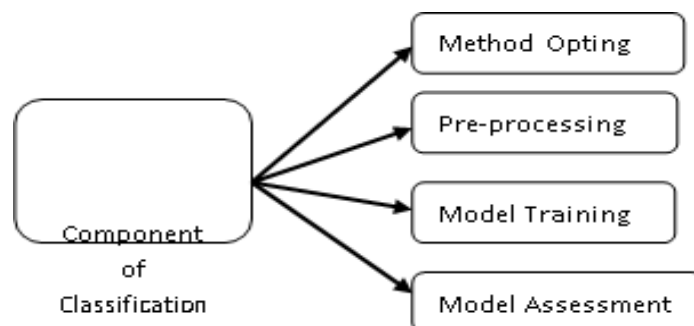


Figure 3.4 Classification Components

METHOD OPTING:

The practices of classifications is impulsive and often moderately nondeterministic, appropriateness of a sole process may not be in promptly acceptable phase. In this situation, the verdict of determining the most admirable forecasting mechanism is vital tricky, predominantly in situation of having modest former awareness about the background of difficulty.

PRE-PROCESSING:

Collated data from various sources may consist some flaws like may encompasses incomplete info set, noise or inconsistency. The pre-processing exercise unite the steps of info cleaning, integration,renovation, lessening and the process of data discretization.

MODEL TRAINING:

The quality of a new build system highly depends on its training procedure. The term training employ to train a system for predict an occurrence of suspicious activity more accurately as event is entered intothe system and confirm that the process data is into an appropriate form for mining. Typically, the new build

system trained with the available tags of the training dataset. In non-availability case of training dataset, the buildmodel adopts the unsupervised technique for training the system.

MODEL ASSESSMENT:

After building a model the assessment procedure plays a significant role in the case of analysing and enhancing the efficiency of a build system. Sometime the build system which produces high resultswith the training dataset may fail to maintain its effectiveness with new data setting. Assessment procedure of building model may solve these issues.

Classification Approaches of Data Mining

Huge of published efforts [6-10], has denoted that from the born age of classification methodology a decent amount of classification algorithms has introduced by investigators for optimizing prediction precision of data with small number of false alarms but every handy practice has associated some unique restraint. Some of the main classification practices of datamining are (fig.2)

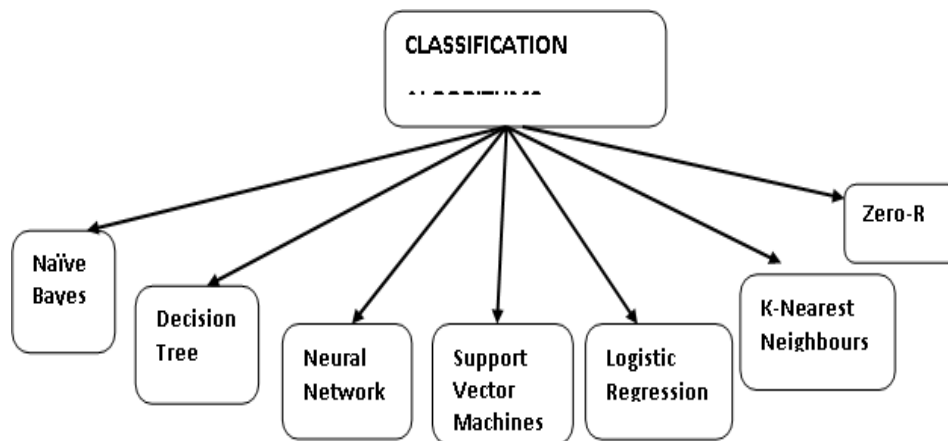


Fig.2 Data Mining Classification Procedures

□ **Naïve Bayes Algorithm:** From several of supervise learning practices this scheme is one most adopted practice for the process of classification [11]. Straightforwardly based on probability theory this scheme utilizes Bayesian theorem for the forecasting practice. Due to its easy constructed model, this procedure is one most fitted scheme into a high contributions environment for manipulating probability through the utilization of following equation.

$$P(v|z)P(z) = \text{-----} \tag{1}$$

$$P(v) \text{ i.e. } P(z|v) = P(v_1|z) \times P(v_2|z) \times \dots \times P(v_n|z) \times P(z) \tag{2}$$

Where :

P(z) & P(v) indicate the class & predictor former prospects. During training phase of model each class possibility is calculated by discovery of the entire occurrence of it into the assorted set of training data, recognized as “prior probability”. The P(z|v) demonstrates posterior opportunity of target class z given attribute predictor v.

□ **Decision Tree:** To settle course of action or show statistical possibility this practice recursively fragments a dilemma into subclasses till it resolve problem straightforwardly without any distress. The problem attributes become into a form of constructed tree nodes and their proportional values regulate tree paths. Due to modest working methodology and strong productivities this practice has led over various approaches.

□ **Neural Network:** In a situation where maximum variables of assorted data consist weakly relevance this scheme outperforms in comparison of other handy algorithms. Major of offered framework of this schemes can be recognized under two name Feed-forward & Recurrent neural networks. The Feed-forward neural network is one of simple network where info explorer into a sole direction only, in to out neuron through theutilization of midway neurons.

□ **Support Vector Machines:** This scheme is a training algorithm to learn classification & regression directions from assorted dataset [8]. This procedure was familiarized in 1960s for the task of classification. It utilizes

the theory of statistical learning for classifying wholly training instances appropriately, separating data under correct relative classes [9]. Mapping of an exclusive nonlinear data into a higher-dimensional feature a hyperplane is raised that bisect classes info and exploit separation margin amid itself and the adjoining points to it. The hyperplane exploited as the foundation to classify unidentified statistics. This practice supports both regression & cataloguing tasks. Furthermore, this practice follows structural risk minimization attitude, closely linked to the theory of regularization. Separately this practice is an appropriate scheme to grip multiple continuous & categorical variables. For categorical variable quantity an imitation data is shaped with values either 0 / 1.

□ **Logistic Regression:** This procedure is one other powerful tool of data modeling, measured as a modest standard statistical style which usages one or more than one predictor (numerical / categorical) for the forecasting tasks. To forecast data of a binary variable, the method of linear regression is not a proper solution.

□ **K Nearest Neighbors Algorithm:** This practice is a non-parametric lazy learning algorithm, not famine any training process and not build assumptions for generalization. In simple form this practice stores wholly accessible cases and categorizes new-fangled cases on the base of similarity measure. Typically, under this practice new arrival instance is evaluated with deposited instances to assign respective class, coming instance linked with relative class through its closest instance. However, over the simple recognition issues this practice executes in a healthy way but in real time practices face many hitches.

□ **ZERO-R:** For data forecasting practice this mechanism depended on the validity of training set. Characteristically, this practice has not supremacy of data forecasting. This practice purely forecast info that fits to the majority category classes. With numeric data this practice forecast an average worth of targeted attributes from training set.

A huge number of other classification practices of data mining is in accessible mode but description of each method is not possible to incorporate in this paper, can be fetched by related published efforts [12-18].

Challenges of Reachable Classification Practices

However, to implement an efficient forecasting method huge effort have puts by related research community over past decades but due to inefficiencies of each offered practices filed is still open for the further research. One foremost failure cause of most offered approaches is their build mechanism, highly depend over the working methodologies of classical procedure. As naïve offered schemes adopted the forecasting functionality of traditional sole scheme hence same faces hitches as traditional mechanism gain during forecasting practice. Furthermore, most of offered scheme not designed for forecasting of different info types, models have built for forecast of only a particular type of info. Such hitches of present schemes show the scope of further research. Separately from such issues some of the other significant restraint can be direct as

- In highly input phase most schemes struggle to maintain their performance.
- Most of handy approach not able to upgrade working functionality without an external effort of expert.
- Offered less precision value with and disadvantage of inflexibility.
- Huge amount of offered scheme fail to forecast naïve info, only recognize info that was used with training data set.

Moreover, implementing cost of most presented forecasting schemes are very high. As related literature of this field indicated that huge research community of this related filed has offered vast practices to forecast info at an acceptable range of accuracy but dissimilar process act differently and still struggles to maintain forecasting act with moderation of info set.

CONCLUSION

Classification practice is a technique to manage data under diverse classes, separated on the base of certain significant characteristics and features. Lot of published efforts has denoted that classification practices are highly useful for machine learning procedures. In modern time frame a huge amount of classification procedures is in accessible mode and rapidly naïve methodologies are proposing by related filed investigators. This paper present some of the accessible classification techniques of data mining. Furthermore, some current challenges of handy algorithms has also elaborated to better explain the modern hitches which may help of naïve investigators to set future work directives.

REFERENCES

- [1] Mahajan, A., & Ganpati, A. (2014). Performance evaluation of rule based classification algorithms. *International Journal of Advanced Research in Computer Engineering & Technology (IJARCET)*, 3(10), 3546-3550.
- [2] Neelamegam, S., & Ramaraj, E. (2013). Classification algorithm in data mining: An overview. *International Journal of P2P Network Trends and Technology (IJPTT)*, 4(8), 369-374.
- [3] Cleary, J. G., & Trigg, L. E. (1995). K*: An instance based learner using an entropic distance measure. In *Machine Learning Proceedings 1995* (pp. 108-114). Morgan Kaufmann
- [4] Umadevi, S. & Marseline, Jeen. (2017). A survey on data mining classification algorithms. Pp.-264-268.

- [5] P. Sharma, S. Saxena and Y. Mohan Sharma, "An Efficient Decision Support Model Based on Ensemble Framework of Data Mining Features Assortment & Classification Process," 2018 3rd International Conference on Communication and Electronics Systems (ICCES), 2018, pp. 487-491, doi: 10.1109/CESYS.2018.8723882.
- [6] Sharma Y.M., Saini P.K., Shalini, Sharma N. (2021) Effective Decision Support Scheme Using Hybrid Supervised Machine Learning Procedure. In: Goyal D., Gupta A.K., Piuri V., Ganzha M., Paprzycki M. (eds) Proceedings of the Second International Conference on Information Management and Machine Intelligence. Lecture Notes in Networks and Systems, vol 166. Springer, Singapore.
- [7] M. Chakraborty, S. K. Biswas and B. Purkayastha, "Data Mining Using Neural Networks in the form of Classification Rules: A Review," 2020 4th International Conference on Computational Intelligence and Networks (CINE), 2020, pp. 1-6, doi: 10.1109/CINE48825.2020.234399.
- [8] Shalini, Saini P.K., Sharma Y.M. (2021) An Intelligent Hybrid Model for Forecasting of Heart and Diabetes Diseases with SMO and ANN. In: Shorif Uddin M., Sharma A., Agarwal K.L., Saraswat M. (eds) Intelligent Energy Management Technologies. Algorithms for Intelligent Systems. Springer, Singapore.
- [9] Wijaya SH, Pamungkas GT, Sulthan MB (2018) Improving classifier performance using particle swarm optimization on heart disease detection. In: IEEE international seminar on application for technology of information and communication (iSemantic), pp 603–608
- [10] Liu X, Wang X, Su Q, Zhang M, Zhu Y, Wang Q, Wang Q (2017) A hybrid classification system for heartdisease diagnosis based on the RFRS method. Hindawi Comput Math Method Med:1–11
- [11] Sharma Shashikant, Soni Vineeta, Pradhan Nitesh (2016), Efficient Technique for Boosting Attack Detection Rate Over a Host or Network System, International Journal of Computer Application, Volume 147 – pp 37-46
- [12] Jaiswal O., Saini P.K., Shalini, Sharma Y.M. (2021) Analyze Classification Act of Data Mining Schemes. In: Goyal D., Gupta A.K., Piuri V., Ganzha M., Paprzycki M. (eds) Proceedings of the Second International Conference on Information Management and Machine Intelligence. Lecture Notes in Networks and Systems, vol 166. Springer, Singapore.
- [13] A. Mehndiratta, D. Singh and Y. M. Sharma, "A new hybrid scheme to improve security for digital message," 2017 International Conference on Inventive Computing and Informatics (ICICI), 2017, pp. 149-153,
- [14] Morteza Nagahi, Raed Jaradat, Mohammad Nagahisarchoghaei, Ghodsieh Ghanbari, Sujun Poudyal, Simon Goerger, "Effect of Individual Differences in Predicting Engineering Students' Performance: A Case of Education for Sustainable Development", Decision Aid Sciences and Application (DASA) 2020 International Conference on, pp. 925-931, 2020.
- [15] Y. Bengio, J. M. Buhmann, M. J. Embrechts and J. M. Zurada, "Introduction to the Special Issue on Neural Networks for Data Mining and Knowledge Discovery", *IEEE Transactions on Neural Networks*, vol. 11, no. 3, pp. 545-549, 2000.
- [16] Verma A, Mehta S (2017) A comparative study of ensemble learning methods for classification in bioinformatics. In: 7th international conference on cloud computing, data science and engineering— confluence. IEEE, pp 155–158
- [17] Bužić D, Dobša J (2018) Lyrics classification using Naive Bayes. In: IEEE 41st international convention on information and communication technology, electronics and microelectronics (MIPRO), pp 1011–1015
- [18] Subitha Sivakumar, 2Sivakumar Venkataraman and 2Asherl Bwatiramba "Classification Algorithm in Predicting the Diabetes in Early Stages" Journal of Computer Science, 2020, 16 (10): 1417.1422.
- [19] Saouabi, M., & Ezzati, A. (2020). Data mining classification algorithms. *Computer Science*, 15(1), 389- 394.