

EFFECTIVE MODEL SELECTION FOR DATA MINING IN HEALTHCARE SERVICES

Swapna Bhavsar

Research Scholar

University of technology, Jaipur

Dr. Ashish Chaurasia,

Professor

University of Technology, Jaipur

DECLARATION: I AS AN AUTHOR OF THIS PAPER/ ARTICLE HERE BY DECLARE THAT THE PAPER SUBMITTED BY ME FOR PUBLICATION IN THIS JOURNAL IS COMPLETELY MY OWN PREPARED PAPER. I HAVE CHECKED MY PAPER THROUGH MY GUIDE /SUPERVISOR /EXPERT AND IF ANY ISSUE REGARDING COPYRIGHT /PATENT /PLAGIARISM/ OTHER REAL AUTHOR ARISE THE PUBLISHER WIL NOT BE LEGALLY RESPONSIBLE IF ANY OF SUCH MATTER OCCUR PUBLISHER MAY REMOVE MY CONTENT FROM THE JOURNAL

ABSTRACT

Data mining is the process of extracting useful information from large stores of information that are typically in an unorganised format. It is put to use in the processing of automated extraction as well as the production of predicted information derived from big records. Data mining technologies have been employed in the medical field for a variety of purposes, including the prediction of prognosis and diagnosis, the assessment of outcomes, the improvement of patient care and decision-making, etc. The purpose of this research study is to investigate the various applications of data mining techniques in the management of public health care at various levels. This will be accomplished by developing data mining-based decision making models that will assist health care administrators as well as planners in determining the opening of First Referral Units or CHC in a more methodical manner and in ensuring that specialised services are provided to all districts in the same manner.

Keywords: Data, Mining, Health, Patient, Public.

1. INTRODUCTION

Data mining refers to the process of using a variety of methods to detect and filter out information or knowledge included in data, as well as extracting this data for use in a range of different fields, such as decision support and prediction, for example. The medical industry as well as the government are able to mine massive amounts of data for useful information on trends, people, and demographic groups thanks to data mining. The process of data mining requires the application of several statistical and mathematical methods, including neural networking, cluster analysis, and predictive modelling. The process of analysing

source data from a variety of perspectives in order to categorise it and turn it into useful knowledge that can be applied by a large number of people in society is known as "data mining." Data mining is a technique that looks for patterns within datasets and creates models from those patterns. The several different sorts of patterns that can be found, the discovery of which is heavily reliant on the data mining tasks that are carried out. Data Mining is gaining popularity in commercial sectors that require the analysis of large amounts of data to unearth knowledge, which can then be transformed into information that is of practical use. These

sectors include: The use of data mining yields fruitful results in situations where the time dimension is applied to the technique to determine the development of the present data and the relations that need to be noticed. The primary reason or objective for which data mining is receiving an increasing amount of attention in the modern period is due to the general availability of a large number of information and the necessity of utilising such a technique in order to convert that information into knowledge. Using various data mining approaches, a significant amount of study work has been carried out on the medical sector all around the world.

2. LITERATURE REVIEW

H. Alaskar (2022) recent ground breaking advancements in deep learning (DL) have been the driving force behind multiple breakthrough results in a variety of difficult computer vision applications. Deep learning neural networks (DNNs) began to outperform shallow machine learning models on a number of key benchmarks in 2012. This marked the beginning of the astounding triumphs and milestones that have been accomplished since then. The field of computer vision has made significant strides forward thanks to the completion of extremely difficult picture interpretation tasks with remarkable precision. These accomplishments have shown considerable promise in a wide number of disciplines, most notably in medical image analysis by paving the way for earlier illness diagnosis and treatment prospects. In recent years, the application of the DNN for object localization has gained the attention of researchers due to its success over conventional methods, particularly in object localization. This is particularly due to the fact that the DNN can localise objects more accurately than conventional methods. This article gives a brief analysis of DNN implementation for medical

imaging and evaluates its performance on benchmarks in light of the fact that this field has grown quite comprehensive and is expanding at a rapid rate. This paper presents the first review that focuses on object localisation using the DNN in medical pictures. The review is focused on this topic because it is presented in this study. This study's primary objective was to provide a summary of recent research that was based on the deep neural network (DNN) for the purpose of medical image localization and to highlight the research gaps that can provide useful ideas for the direction of future research concerning object localization tasks. To begin, an overview of the significance of medical image analysis and the many technologies now available in this field is provided. After that, we will talk about the DNN that is most commonly used in the recent research that has been published. Finally, we bring this article to a close by examining the difficulties connected with the application of the DNN for medical picture localization. These difficulties can be the impetus for more research that identifies potential future advancements in the applicable field of research.

Lijing Ren (2022) As a result of the intelligent connections made possible by the Internet of Medical Things (IoMT), the quality and output of the medical business have significantly increased. Protecting the privacy of patient information in the face of emerging challenges, such as the transfer of biological data through open and untrusted networks, is essential. The inability of IoMT to process data effectively prevented it from making full use of more conventional forms of encryption to safeguard sensitive information. In this article, we present a unique data protection model that we built specifically for medical photos. The concept makes use of visual cryptography (VC) to store biomedical data in a separate database. This allows the sensitive data of patients to be transferred in a straightforward and risk-free

manner. We continue to make use of transfer learning in order to train an improved neural network in order to mitigate the decrease in biological identification performance that is brought on by VC-based noise. The findings of the experiments indicate that the suggested strategy protects users' privacy within the IoMT environment while also preserving the high level of accuracy achieved by biomedical recognition.

Judith Santos-Pereira (2022) the field of healthcare has gotten more difficult in recent years, necessitating the extraction of useful information from vast quantities of intricate data in order to locate the most effective therapies. Although a number of studies have proposed utilising Data Mining techniques as a solution to the problems, none of these works have indicated a particular tool as being the most effective solution. This paper presents a survey of popular open-source data mining tools in which data mining tool selection criteria based on healthcare application requirements is proposed and the best ones using the proposed selection criteria are identified. The purpose of this paper is to fill this gap by providing an overview of popular open-source data mining tools. KNIME, R, RapidMiner, Scikit-learn, and Spark are some of the well-known open-source data mining tools that are evaluated in this article. According to the findings of the study, KNIME and RapidMiner offer the most comprehensive coverage of the requirements for healthcare data mining.

Ahmed Mahdi Abdulkadium et al (2022) While advances in technology, particularly in the form of computer-based healthcare information apps and hardware, are making it easier to collect data about healthcare and gain access to that data, these advancements are still limited. In this scenario, there are technologies that can be used to analyse and investigate this

medical data once it has been obtained and stored. An examination of the documented medical data records may assist in the identification of concealed characteristics and patterns, which could considerably expand our comprehension of the beginning of sickness and the treatment therapies. Significantly, the development of information and communications technology (ICT) has surpassed our capacity to evaluate, summarise, and glean insight from the data. This is a problem. The development of database management systems has provided us with the fundamental tools necessary for the efficient storage as well as lookup of massive data sets; however, the question of how to make it possible for humans to interpret and analyse large amounts of data is still a difficult and unresolved problem. Therefore, in order to deal with vast amounts of data, advanced methods for automated data mining and the discovery of new information are required. Within the context of this investigation, an effort was made to obtain knowledge that will assist various individuals in making judgments that will guarantee that the sustainability targets on Health are realised by utilising a method known as machine learning. In conclusion, the most recent data mining methodology, together with data mining methods and also its deployment tools, which are more beneficial for healthcare services, are discussed in length.

O V Klochko (2022) this article discusses the features of applying regression analysis methods in machine learning systems. Data mining is presented in this article. Data mining is based on the methodologies of mathematical statistics and machine learning. The machine learning model that was developed contains regression analysis modules that are based on linear regressions, Bayesian linear regressions, artificial neural networks, decision trees, and decision forests. The procedure of using this machine learning model included employing

the techniques that were discussed, constructing the matching regression models, doing a comparative analysis of the models, and analysing the outcomes. The findings that were collected provide evidence that it is possible to use data mining in medical research by utilising machine learning algorithms. The methodologies that have been given have the potential to serve as a foundation for the strategic development of new avenues for the processing of medical data and decision-making in this area. We have recognised the opportunities for additional study that are available in the field of applying data mining techniques to the healthcare system, specifically clustering, classification, and anomaly detection.

3. RESEARCH METHODOLOGY

In order to facilitate data mining in the decision making process for public healthcare management analysis, the primary concern in the process of building conceptual models must be addressed. The model demonstrates how decision-making processes can benefit from the application of data mining tools. The CRISP-DM (CRoss-Industry Standard Process for Data Mining) Reference model served as the inspiration for the design of the model that was proposed. CRISP-DM is a data mining process model that was established by the leaders of the industry in partnership with experienced data mining users and data mining software tool vendors. CRISP-DM was given the name "CRISP" after the acronym for the CRISP Data

Mining Process. In the body of research that has been done, further data mining process models have been found. They are conceptually comparable to CRISP-DM, despite the fact that they make use of slightly different language. Following examination during the data mining process, the CRISP-DM model was found to be the one that was best suitable for the proposed framework. The CRISP-DM reference model divides the activities of data mining into the six phases listed below. Each of these phases includes a number of tasks, including business understanding, data understanding, data preparation, modelling, assessment, and deployment. In order to meet the requirements of the proposed model, the CRISP-DM reference model has been modified.

4. DATA ANALYSIS

- At 37 percentile (as indicated by the brown line), the gain% is 93.57, which shows an excellent performance of the decision tree created by employing the C 5.0 Algorithm. This information can be found in the gain chart of the national Model, which can be found in Fig.1.
- When compared with both the actual outcome and the predicted outcome, the model has an accuracy of 95.99%, which is sufficient to satisfy the criteria established by the research.

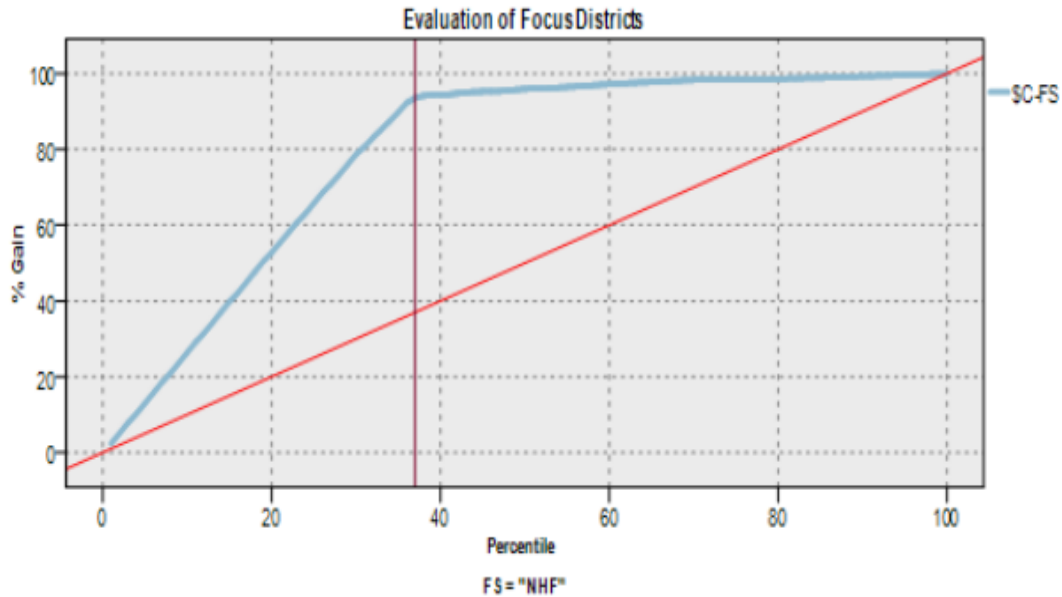


Figure 1: Gain chart for National model

Coincidence Matrix:

The Coincidence Matrix displays, for categorical targets, the pattern of matches that may be found between each generated (predicted) field and its respective target field (flag, nominal, or ordinal). When there are target categories, it assists in making more accurate predictions. A table is presented with rows that are defined by real values and columns that are defined by predicted values.

Within each cell of the table is the number of records that match a particular pattern. According to the coincidence matrix that is displayed in table 1, there are a total of 420 focus state districts, and 410 of those districts meet the requirements to be classified as High Focus Districts. 237 of the non-focus districts out of a total of 254 meet the requirements to be considered non-focus districts. The effectiveness of the coincidence matrix is fairly satisfactory.

Table 1: Coincidence Matrix For \$C-FS (Rows Show Actual) National Model

	HF	NHF
HF	410	10
NHF	17	237

Evaluation Metrics: The area under the curve (AUC) and the Gini coefficient are the evaluation metrics that are produced while evaluating binary classifiers. These two

evaluation metrics, along with their respective scores, are derived concurrently for each binary model. The values of the metrics that are presented in show that AUC and GINI values

are relatively high, which demonstrates that the model achieved a good level of accuracy.

Model	AUC	Gini ‘
\$C-FS	0.972	0.945

Variables ranking validation

Owing to the fact that our goal is to identify the model that has a lower total number of variables and a greater percentage of correct answers. As

can be seen in table 2, the ranking of the variables in this model provides an indication of the primary healthcare variables that are involved in the usage of healthcare services.

Table 2: Attribute/Variable Ranking Validation

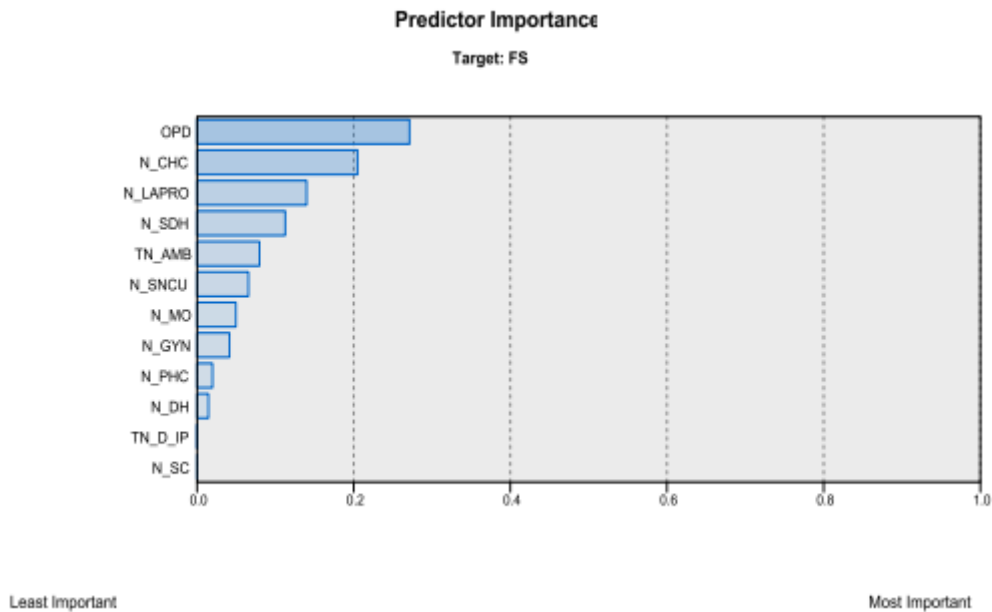
Ranking	PI using data of 2013-14	PI using data 2012-13
One	OPD	OPD
Two	N_CHC	N_CHC
Three	N_LAPRO	N_SDH
Four	N_SDH	TN_AMB
Five	TN_AMB	N_SNCU
Six	N_SNCU	N_LAPRO
Seven	N_MO	N_MO
Eight	N_GYN	TN_D_IP
Nine	N_PHC	N_SC
Ten	N_DH	N_DH

It is quite encouraging to see that the accuracy of the model is greater than 90% for both the 2013-14 and 2012-13 academic years. The rankings of the variables are the same in both years, despite the fact that the variables OPD and N CHC are in the same position. This

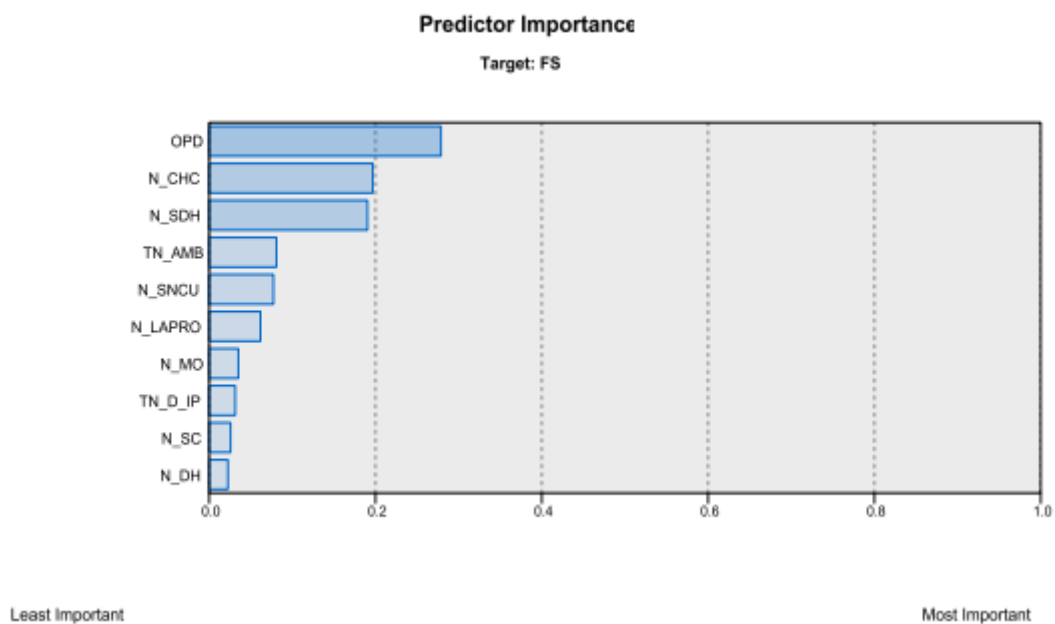
demonstrates that the combinations of variables play an extremely crucial role in the process of estimating the availability of a variety of healthcare services. Both of the years' segments contain the variables OPD, N CHC, N SDH, N SNCU, N LAPRO, and N MO. The TN AMB

variable also appears in both of the years' segments. We are easily able to make accurate projections about the requirements for particular medical facilities in the districts that are in question thanks to the small number of

factors. In addition, these are highly beneficial in terms of time-series analysis and forecasting because they have an accuracy of greater than 90%.



Predictor Importance Chart using data of 2013-14



Predictor Importance Chart using data of 2012-13

Figure 3: Comparison of Predictor Importance Charts

5. CONCLUSION

The results of both analyses of the data suggest that the OPD, N CHC, N SDH, T AMB, and N SNCU variables are the most relevant ones. This demonstrates that the priority does not change regardless of which focus districts are being considered. The Community Health Center (CHC) and the State Department of Health (SDH) are both referral units for lower health institutions and serve as a catalyst for increasing delivery services. Because the Non-Focus states' community health centres (CHC) and state district hospitals (SDH) are so well run, they are better able to meet their goal of providing effective care to patients in outpatient and inpatient settings alike.

REFERENCES

1. H. Alaskar (2022),” Deep Learning Approaches for Automatic Localization in Medical Images”,
2. Lijing Ren(2022),” A New Data Model for the Privacy Protection of Medical Images”,
3. Judith Santos-Pereira (2022),”Top data mining tools for the healthcare industry”,
4. Ahmed Mahdi Abdulkadium et al (2022),”Application of Data Mining and Knowledge Discovery in Medical Databases”,
5. V Klochko (2022) ,” Data mining of the healthcare system based on the machine learning model developed in the Microsoft azure machine learning studio “,
6. Mounir El Khatib (2022) ,” Digital Disruption and Big Data in Healthcare - Opportunities and Challenges”,
7. Zafrul Hasan (2021) ,” An Analysis on Data Mining Applications in Healthcare Sector”,
8. Judith Santos-Pereira (2021),” Top Data Mining Tools for the Healthcare Industry “.
9. Xunjie Gou (2021),”An overview of Big Data in Healthcare: multiple angle analyses”,
10. SUNITA VERMA (2021),” Data Mining: An Effective Tool in Cognitive Psychological Emotion Analysis”,
11. Sri Venkat Gunturi Subrahmanya(2021) ,” The role of data science in healthcare advancements: applications, benefits, and future prospects”,