

## A SURVEY ON DATA MINING TECHNIQUES IN HEALTHCARE

V. Krishnakumar<sup>1</sup>, Dr. V. Sangeetha<sup>2</sup>

<sup>1</sup>Research scholar, Department of Computer Science, Periyar University, Salem, India

<sup>2</sup>Department of Computer Science, Periyar University Constituent College of Arts and  
Science, Pappireddipatti, Dharmapuri, India

kichuveera@gmail.com, sangee759@gmail.com

**Abstract** - There are some models available in the field of healthcare to diagnose or predict the dengue fever, early stages. Process of these model to diagnose or predict the dengue fever is difficult and met major drawbacks such as less accuracy, inadequate performance, time, and cost factor is too high also absence of advanced technique and tools. Predictive model have been created to solve drawbacks of previous models. This predictive model come up with data mining, techniques, algorithm, tools and big data. There are some predictive model introduced in the field of healthcare using classification technique but inadequate performance in both diagnosis and predict of dengue fever. This survey paper is organized based on data mining techniques, algorithms and tools built on the list of attribute. The attributes are technique, algorithm, type, approach, interface, IDE, data set collected from, data set collected year, format of the data set, data set size, method for removing missing values, predictive model technique, time complexity, tool, pros, and future works. Based on the existing works, results are dissimilar, accuracy level and outcome of the diagnosis or prediction level is compete each other to prove its better performance. If the appropriate techniques and algorithms are used will improved the performance as well as accuracy levels. This paper gives the summary of the data mining concepts, techniques, algorithms, tools, challenges and discusses future work of data mining techniques.

**Keywords** –Data Mining, Techniques, Algorithm, Predictive model, Tools

### I. INTRODUCTION

Data and information have play vital role as a resources formost of the organizations [1] [2]. In early days data produced mechanism are very less, with the minimum amount of data processing speed is very high. Particular data after enter into process level gives output as information and store that particular information in any one of the storage devices using data base or data warehouse technology. Storage space was adequate, so each and every industry

start to use technology forgetting better result. After few years data produced mechanism are increased such as internet, social media, and mobile data and others. Using these mechanism data have been growing ultra-fast manner. In this scenario storage space was inadequate and get better knowledge from data was difficult. In recent year's industries like educational organizations, bank/finance based detail, healthcare sectors, working group of crime detections and all others are collecting more data continuously than human expected limit. Especially in the field of healthcare, health related data which means, the patient information (including patient id, patient name and complete detail of individual patient), symptoms, prescription details, function of the human part like, heart, kidney and all other parts as media file and issues in healthcare are growing faster than technologies. From the huge collection of data or information could not found better result using existing technologies. In the view of previous technology have lots of drawbacks such as, storage space, size of the data, types of the data and format of data, accuracy, time and cost being present. Data mining come up with solution with its techniques and algorithm for previous technology drawbacks. Big data come to play to solve storage related issues. This paper confers the idea of different types of technique and algorithms and also describes as a report of survey data mining techniques and algorithms. This paper is structured as follows. The related work discussed in section 2 and section 3 gives the process of data mining concepts and section 4 describes data mining in healthcare for dengue and section 5 discussed classification technique used in dengue prediction. The analysis briefed in section 6 and section 7 provides the conclusion and future work.

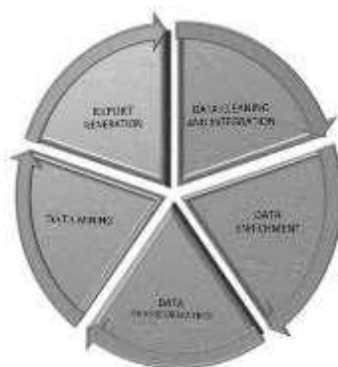
## II. RELATED WORK

Data mining with analysis and analytics in the field of healthcare is very fast growing. The existing research works and various technology backgrounds done in the data mining with healthcare are abridged in this section. This section gives an overview of data mining techniques and algorithm used by the following authors. D. Usha rani surveyed the data mining tools and techniques in medical field in 2017 and deliberated about data mining tools and techniques to be improved for reduce cost and time from human resources and capability. Also, use the combination of data mining techniques than particular technique applied for a specific diseases diagnosing or predicting diseases in healthcare sector could produce more advance level of positive results [3]. A. Shameem Fathima *et al* reviewed knowledge discovery, data mining, and the purpose of the classification methods to diagnosis and prognosis used for

arbovirus-dengue in 2011. The proposed work was apply hybrid classification schemes and implement data mining tools to analyze the data, evaluate the data mining algorithm through these steps try to provide some perceptible health information mined by the data mining methods [4]. Subhash Chandra Pandey explored various tasks of data mining, benefits and drawbacks of data mining techniques in the field of healthcare and compare the health related data, also suggested to get the better result, need to improve the data mining techniques in 2016 [5]. Dr. A. R. PonPeriasamy and S. Mohan summarized the significance of data mining concepts applied in medical healthcare sectors, importance of preprocess task and analyzed the exclusivity of medical data mining in 2017 [6]. Parvez Ahmad *et al* précised different types of data mining techniques, applications and challenges in health care perspective in 2015. Above all discussed the feature selection methods which is one of classification techniques in data mining and importance of secure distributed healthcare environment [7]. R. Naveen Kumar, M. Anand Kumar conceded the preprocess techniques produced a complete data set instead of using raw dataset which leads to delay in the overall process [8]. S. Sharath *et al* surveyed clinical data mining concepts to mine the health related records. Innovation of data mining technology is distributed network which is used to share storage space with legitimate channel [9].

### III. DATA MINING

Data mining used to extract useful knowledge using huge volume of data from past and present activities (collectively called 'large dataset') to give better outcome. Data mining consist different kinds of techniques and algorithm available to diagnose or predict the particular diseases. Classification is one of the best technique in data mining to give accurate result in the field of healthcare using some statistical languages.

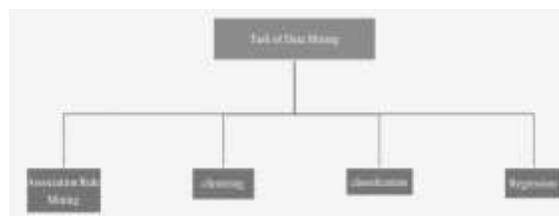


**Fig 1: KDD Process Cycle**

Description of KDD process as follows

1. Data cleaning and integration is the first in KDD process step to remove the no longer used data irrelevant or unwanted data from existing resources. Integration which means combine different format data as a same format.
2. Data enrichment is the second step in KDD process to increase value of data with existing one.
3. Data transformation is the third step in KDD process to put the data in particular format that user want. This is also called format uniformities.
4. Data mining is the fourth step in KDD process to extract the knowledge from huge data repository.
5. Report generation is the fifth step in KDD process to get expected outcome.

#### IV. DATA MINING IN HEALTHCARE FOR DENGUE



**Fig 2: Data mining algorithms used for dengue prediction**

1. Association rule mining work based on relationship between the data points. For example Hospitals collect the symptoms of the patient's data after collect the symptoms can apply this ARM, the hospital can find the number of patient affected by same kind of disease furthermore find the particular diseases spread in which place.
2. Clustering Analysis also called as cluster analysis. For example identify group of patients with same properties.
3. Classification is used to classify the patient by symptoms. For example, patient age, masculinity, blood pressure, presence or absence of certain symptoms, etc.
4. Regression also used in predict analysis based numeric values.

## V. IMPORTANCE OF CLASSIFICATION ALGORITHMS.

One of the best algorithm among data mining algorithms to extract the exact well known category to expound the feature of the model from dataset. This classification algorithm has two phases such as i) Training and ii) Test phases. Evaluation has five constraint factors in classification is i) accuracy ratio of the prediction from the dataset ii) time and cost iii) robustness of all about dataset iv) model to determine of the data and v) understandability of particular model.

Different types of algorithm available with classification type based algorithm like, i) DT algorithm, ii) Rule-Based algorithm iii) statistical algorithm iv) Non-Linear algorithm and so on.

A. DT Algorithm is also called as decision tree algorithm. It consist two important algorithm such as i) iteratedichotomies 3 (ID3) and ii) c4.5(C4.5- Extension of I D3 algorithm)

B. Rule based algorithm – CN2 and CL2

C. Statistical algorithm – cart, genetic algorithm

D. Nonlinear algorithm – neural network, nonlinear regression

## VI. SURVEY ON DATA MINING TECHNIQUES

Different authors' to diagnose or predict the dengue fever using data mining techniques is analyzed and discussed in this section. This survey is based on the following list of parameters. The parameters are technique, algorithm, type, approach, interface, IDE, dataset, method for rmv, pmt, tc, tool, pros, future work. The above lists of parameters are shown in the Table I.

P. Manivannan, *et al* developed algorithm for dengue fever prediction using on k-medoid clustering based on unsupervised algorithm. In this research also used dwin's method for remove missing values. Time complexity calculated as  $O(n^2)$  and practically proved. Data set collected from Seremban District Health Office, Negeri Sembilan, Malaysia in the year of 2010-2013 size of this data was 171 attribute and 1910 records. Proposed work improve evolutionary patterns and apply dengue fever prediction [10].

HUSAM, I.S *et al* created the model using statistical and expert verification tools. Predictive model techniques such as J48, DTNB and naïve bayes are applied in this model. Wrapper approach applied for evaluate the future subsets and PSO, GA, RS algorithms also used. Compared with each other and found result. Dataset obtained from Public Health Department, Seremban, Negeri Sembilan, Malaysia, format of the was excel, year of the data 2003-2010, size the dataset is 20 attribute and 6082 records [11].

Kamran Shaukat *et al* implemented predictive model using classification techniques and WEKA tool. Compared DT, NB and J48 with other. Algorithm such as, naive bay, J48, SMO, REP tree, and random tree also the purpose of interface point of view in this research Explorer, knowledge flow and Experimenter interfaces are used. Dataset took from dataset was collected from District Headquarter Hospital (DHQ) Jhelum. Format of the dataset was in comma separate value (.csv) it consists 18 attribute and 99 records [12].

Ashwini Rajendra Kulkarni *et al* created association rule generation for virus-related illness using association rule mining technique, algorithm like apriori, FP – growth, association rule generation and rapid miner text mining tool. This paper defined abstract level of association rule mining and apriori algorithm and proved through implementation level. Dataset format was in excel. Future work of this paper is produce the numerous patterns [13].

Ramandeep Kaur *et al* explained and evaluated dengue fever in high population area using cluster technique with R-Studio tool for statistical report based on unsupervised learning. User created dataset for this research with .csv file format it consists of 10 attribute and 100 records [14]. Nandini. V *et al* developed a system for detection and prediction of dengue fever using frequency analysis, classification and regression technique, interface designed as patient and researcher GUI and supervised learning type. Various predictive model technique and tools also used such as time series, exponential smoothing, moving average, simple linear regression and tools like SAS and LibSVM. Format of dataset is .csv in the year of 2010- 2015. This dataset consists of 100 records. Future work is create desktop application and browser plugin file [15].

Shaufia *et al* discovered algorithm for identifying DHF and TF using association rule mining technique with apriori, FPGrowth and Intersection Set Theory-Expand FP algorithm (ISTEFP). IST-EFP is also used to remove missing values from the particular dataset. Suggested for future work is reduce the dimension of the dataset using Intersection

Set Theory-ExpandFP algorithm [16].M.V.Jagannatha Reddy *et al* developed expert system usingneural network technique, decision tree algorithm, MATLAB.

2013a tool. Dataset are obtained from Srinagarindra hospitaland ongklanagarind hospital, Thailand and remove the missingvalue manually. Future work taken from this paper is predictany types of fever [17].

M Krishna Satya Varma *et al* developed decision tree modelusing decision tree technique, ID3 algorithm based onunsupervised learning [18].

Dr. ArunKumar.P.M *et al* filtered related to dengue fever data using decision tree, support vector machine techniques as well as its algorithm, Netbeans acted as a user interface, tools forthis work was WEKA and predictive model technique calledfisher filtering and prediction. Important note in this work wasthis model worked in google application, cloud computing andDT. Future work will predict the environment factors [19].

## VIII. CONCLUSION

Understanding of classification technique and developedpredictive model is depends on the dengue related data.Researcher has given overview of different predictive modelusing classification methods with their advantages anddisadvantages. Many approaches differ in the way of comparethe accuracy. Several approaches suggest additional techniquesand also advanced tools for diagnosis or prediction of denguefever. Preprocessing also suggested. Survey revealed thatconsidering classification methods choice is important,because, number of available predictive model provides samekind of functionality. In diagnosis or prediction of dengue feverstatistical language also involved. If the classification, preprocess, cloud computing storage purpose will help thediagnosis or prediction of dengue fever more efficient. The researcher will use data mining completely in predictivemodel classification that will provide tool to diagnose orprediction the dengue fever. Complexity of both time and spaceis not concentrate as much it is also considered for the futureresearch work.

## REFERENCES

- 1.Klosgen W and Zytchow J M (eds.), “Handbook of data mining andknowledge discovery”, OUP, Oxford, 2002.
2. Kantardzic M, “Data mining: concepts, models, methods, andalgorithms”, John Wiley, New Jersey, 2003.

3. D.Usha Rani “A Survey on Data Mining Tools and Techniques in Medical Field ”, International Journal of Advanced Networking & Applications (IJANA) Volume: 08, Issue: 05 Pages: 51-54 (2017) Special Issue.
4. A. Shameem Fathima, D. Manimegalai and Nisar Hundewale “A Review of Data Mining Classification Techniques Applied for Diagnosis and Prognosis of the Arbovirus-Dengue”, IJCSI International Journal of Computer Science Issues, Vol. 8, Issue 6, No 3, November 2011 ISSN(Online): 1694-0814.
5. Subhash Chandra Pandey “Data Mining Techniques for Medical Data: A Review”, International conference on Signal Processing, Communication, Power and Embedded System (SCOPE)-2016.
6. Dr. A. R. PonPeriasamy , S. Mohan “A Review on Health Data Using Data Mining Techniques”, International Journal of Advanced Research in Computer Science and Software Engineering Volume 7, Issue 3, March 2017.
7. Parvez Ahmad, Saqib Qamar, Syed Qasim Afser Rizvi “Techniques of Data Mining In Healthcare: A Review”, International Journal of Computer Applications (0975 – 8887) Volume 120 – No.15, June 2015
8. R. Naveen Kumar, M. Anand Kumar “Medical Data Mining Techniques for Health Care Systems”, International Journal of Engineering Science and Computing, April 2016.
9. S. Sharath, M. N. Rao, H. G. Chetan “A Survey on the principles of mining Clinical Datasets by utilizing Data mining technique”, International Journal of Innovative Research in Computer and Communication Engineering (An ISO 3297: 2007 Certified Organization) Vol. 2, Issue 4, April 2014.
10. Kamran Shaukat, Nayyer Masood, Sundas Mehreen and Ulya Azmeen, “ Dengue Fever Prediction: A Data Mining Problem “, Data Mining in Genomics & Proteomics (Dar et al., J Data Mining Genomics Proteomics) 2015, 6:3.
11. M.V.Jagannatha Reddy and B.Kavitha, “ Expert System to Predict the Type of Fever Using Data Mining Techniques on Medical Databases” , International Journal of Computer Sciences and Engineering Volume-03, Issue-09 2015.



12. M Krishna Satya Varma, N K KameswaraRao, "Dengue data analysis using decision tree model", International Conference On Emerging Trends in Science Technology Engineering and Management 09th & 10th, October 2015.
13. Nandini. V, and Sriranjitha. R and Yazhini. T. P, "Dengue Detection and Prediction System using Data Mining with Frequency Analysis", Natarajan Meghanathan et al. (Eds): ACITY, VLSI, AIAA, CNDC pp. 53–67, 2016. © Computer Science & Information Technology (CS & IT)- CSCP 2016.
14. Shaufiah and Bobby Siswanto, "Association Rule Mining For Identifying Dengue Hemorrhagic Fever (DHF) and Typhoid Fever (TF) Disease with IST-EFP Algorithm", Fourth International Conference on Information and Communication Technologies (ICoICT) 2016 IEEE.
15. P. Manivannan, Dr. P. Isakki @ Devi, "Dengue Fever Prediction using K-Medoid Clustering Algorithm", International Journal of Innovative Research in Computer and Communication Engineering, Vol. 5, Special Issue 1, March 2017.
16. Husam, I.S., abuhamad, azuraliza abubakar, suhailazainudin, Mazrurasahani & zainudin mohdali "Feature selection algorithms for Malaysian dengue outbreak detection model (Pemilihan Ciri Algoritma untuk Model Pengesanan Wabak Denggi)", Sains Malaysiana 46(2)(2017): 255–265.
17. Ashwini Rajendra Kulkarni, Dr. Shivaji D. Mundhe, "Data Mining Technique: An Implementation of Association Rule Mining in Healthcare", International Advanced Research Journal in Science, Engineering and Technology Vol. 4, Issue 7, July 2017.
18. Ramandeep Kaur, Gaurav Gupta, "Demographic Analysis of Dengue Fever Using Data Mining", International Journal of Advanced Research in Computer Science Volume 8, No. 7, July – August 2017.
19. Dr. Arun Kumar. P.M, Chitra Devi. B, Karthick. P, Ganesan. M and Madhan. A.S, "Dengue Disease Prediction Using Decision Tree and Support Vector Machine", "SSRG International Journal of Computer Science and Engineering- (ICET'17) - Special Issue - March 2017.
20. M. Bhavani, and S. Vinodkumar, "A data mining Approach For Precise Diagnosis Of Dengue Fever", International Journal of Latest Trends in Engineering and Technology Vol. (7) Issue (4), pp. 352-359.

21. Chia-Bao Chu, Chao-Chun Yang, “Dengue-associated Telogen Effluvium: A report of 14 patients”, *dermatologica sinica* 35 124-126 Mar30, 2017.
22. Buchade Omkar, Dalsania Preet, Deshpande Swarada, Doddamani Poonam,” Dengue Fever Classification using SMO Optimization Algorithm”, *International Research Journal of Engineering and Technology (IRJET)* Volume: 04 Issue: 10 |Oct -2017.
23. Jastini Mohd Jamil, Izwan Nizal Mohd Shaharane and Ve Chun Yung, “An Innovative Data Mining and Dashboard System for Monitoring of Malaysian Dengue Trends”, *Journal of Telecommunication, Electronic and Computer Engineering* Vol. 8 No. 10.



S. No	Technique	Algorithm	Type/ Approach/ Interface/ IDE	Dataset					Method For RMV	PMT	TC	Tool	PROS	Future Work
				DSCF	Y	F	SDS							
							A	R						
1	K-Medoid [10]	Clustering	Unsupervised Algorithm	HCMC	2010 - 2013	-	171	1910	Dwin's	-	O(n <sup>2</sup> )	-	Proof Based	Enhance evolutionary Patterns and apply hierarchical clustering Algorithms for dengue fever prediction.
2	Feature Selection [11]	PSO, GA, RS	Wrapper Approach	PHD, Seremban, neger seremban,	2003 - 2010	Exce 1	20	6082	-	J48, DTNB And Navie bayes	-	Statistica l and expert verificati on	Comparison	-

				Malaysia										
3.	Classification [12]	Naïve Bays, J48 tree, SMO, REP tree & Random Tree algorithms	Explorer, knowledge flow and Experiment er interfaces	(DHQ) Jhelum	-	.csv	18	18	-	-	-	WEKA	Comparison of DT,NB and J48	-
4	Classification [13]	REP Tree, J48, SMO, ZeroR and Random Tree	-	Created	-	.csv	10		-	-	-	WEKA	i)High accuracy ii)Compared by plotting graphs and table	-
5	ARM[14]	Apriori,FPgrowth and association rule generation algorithm	-	-	-	Excel	-	-	-	-	-	Rapid Miner (Text Mining)	-	Apply text mining to produce the various patterns from healthcare data.
6	Clustering [15]	-	Unsupervised Learning	Created	-	csv	10	100	-	-	-	R - Studio	-	-
7	Classification [16]	Frequency Analysis, Classification	Supervised Learning/ Patient and		2010 - 2015	.CSV		100	-	Time series, Exponentia	-	SAS (use logistic	UMLS	Implement desktop app and

		Regression	Researcher GUI							Linear Smoothing, moving Average, Simple Linear Regression		regression model) LibSVM		Browser plugin.
8	ARM[17]	Apriori and FP-growth		Hospital Medical Record	-	-	9	192	IST-EFP	-	-	-	Dimensional reduction. Predictive DHF and TF.	reduce dimension of the dataset using ISTEFP algorithm
9	Neural Network[18]	Decision Tree Algorithm	Decision Tree	Srinagarindra Hospital and Songklanagarind Hospital, Thailand	-	Categorical to Numerical data	>400	-	Manual missing value imputation	Multivariate Model	-	MATLAB 2013a	i)generate performance curves, ROC curves, Confusion curves for both training and test data. ii)an accuracy of 100.0% in children and adults using both	Extended to predict any type of diseases.

													clinical and laboratory features	
10	Decision Tree [19]	ID3	Unsupervised Learning	Health department, Hospital, Urban Local Body	-	-	6	142	-	-	-	-	supervised classifier model	-
11	Decision Tree,SVM and ANN[20]	SVM	Net Beans	-	-	-	18	108	-	Fisher Filtering and Prediction.	-	WEKA Feature Selection	Google Applications Cloud Computing DTS	Environmental factor prediction
12		Optimization [18]	Probabilistic Neural Network	GDS5093	-	-		56	-	-	-	Greedy forward selection algorithm	Implementation SMO	Non – Dominated sorting genetic algorithm – II

**Table I - Analysis of various techniques and tools used in data mining for dengue fever**