# Clinical Decision Support System for Heart Diseases Using Extended Sub Tree

Kulkarni Ravindranath Bhaskar[1], Kulkarni Radhika Ravindranath[2], Kulkarni Rashmi Ravindranath[3]

[1]BE Mechanical, Vaidya College, Pune, Maharashtra, India
[2]BE ENTC, MKSSS's CCOE, Pune, Maharashtra, India
[3]ME Computer, MCOE, Pune, Maharashtra, India

*Abstract*- **Clinical Decision Support System (CDSS) is a tool which helps doctors to make better and uniform decisions. There are many existing systems present which are used for diagnosing the diseases. For different types of diseases the existing CDSS systems changes with different algorithmic approaches. Every approach has its pros and cons. Selecting the positive aspect and overcoming the problems is the main motive.**

**This paper focuses on comparative study of existing CDSS systems namely Mycin, DeDombal, Quick Medical Record (QMR), Internist 1. Also the paper focuses on different algorithmic approaches for CDSS. It also give comparative study for algorithmic approaches of heart diseases. The proposed system deals with the similarity matching function. In decision tree construction, the nodes are constructed on splitting attribute or the flag value. Hence if continuous value is to be handled then it can prove fatal. This kind of flaw is observed in ID3 algorithm. Hence to overcome this drawback Extended sub tree approach is implemented. The results show the comparative analysis between these two approaches.**

*Index Terms*- **CDSS, Patient health Information, Electronic Medical Record, Extended Sub tree (EST).**

## I. INTRODUCTION

Clinical Decision Support (CDS) systems provides clinicians, staff, patients and other individuals with knowledge and person specific information , intelligently filtered and presented at appropriate times, to enhance health and healthcare[2]. CDSS is a tool to assist user in taking clinical decisions of diagnosis. A typical user of CDSS is a physician, nurse or any other paramedical service provider. It gathers the patient health information (PHI) entered by the user in the system. Using pre-determined algorithms or rules, CDSS provides clinically relevant information and conclusions to the user. The rules used in the system can be configured by the administrator. Security of each patient's personal record must be provided[1].

## II. DIFFERENT EXISTING SYSTEMS

Different CDS Systems that were developed from the early times have brought up to professional's attention in 1950's. De Dombal's system was developed at university of Leeds in the early 1970's by deDombals and his associates. They studied the diagnoses process and developed a computer-based decision aids using Bayesian probability theory [Musen, 2001]. INTERNIST-I was a broad-based computer-assisted diagnostic tool developed in the early 1970's at the University of Pittsburgh as an educational experiment [Miller et al., 1982; Pople, 1982]. MYCIN was a rule-based expert system designed to diagnose and recommend treatment for certain blood infections (antimicrobial selection for patients with bacteremia or meningitis) [Shortliffe, 1976].

Table 1: Existing Systems

| Sr No. | Properties | MYCIN | De Dombal | Internist-1 | DXplain | Quick Medical Reference (QMR) |
|---|---|---|---|---|---|---|
| 1. | Developed By | Stanford University | University of Leeds | University of Pittsburgh | Laboratory of Massachusetts General Hospital | University of Pittsburgh |
| 2. | Year | 1970 | 1972 | 1970 | 1970 | 1970 |
| 3. | Diseases | blood infections | abdominal pain | knee replacement surgery | 2,200 unique diseases | Abdomen Pain Severe, Blood Hepatitis |
| 4. | Classification Approach | IF-THEN rules | Bayesian probability theory | Bayesian probability theory, Decision Tree | probabilistic algorithm | Basic Decision Tree |

## III. HEART DISEASES

Heart is the vital organ of the body. Without heart the living organism cannot survive. The working of the heart is only to pump the blood in and out. This creates blood circulation in entire body. Blood circulation helps other organs to work efficiently into the body. There are no.of factors which affect heart to malfunction such as history of patient as well as hereditary , life style , poor diet, high blood pressure, obesity, percentage of cholesterol, high per tension, smoking and drugs habbits etc[7].

## IV. DIFFERENT APPROACHES FOR DIAGNOSING HEART DISEASES

There is large amount of heart related data present, which is in unstructured format. Hence by analyzing the data and formatting it into structured manner helps for making the decision. For diagnosing the disease there are many ways in which heart related diseases can be diagnosed and treatment can be provided.

Different approaches have different aspects in diagnosing the diseases.  By using the Neural network approach the accuracy secured was around 80- 90% but the hidden layers description cannot be evaluated [5]. In fuzzy logic approach the weighted rules are generated initially and then the fuzzy rule decision is provided [5][6] and the accuracy obtained id around 79.05%. In naive bayes classification approach helps in predicting whether the patient is prone to heart disease or not and depicting the risk factor for heart attack [7]. The accuracy observed for naive bayes approach was around 90% [8]. Similarly by using Support vector machines concept the accuracy was achieved around 84.12%. While as by using decision tree approach the accuracy increased up to 96% [8].

Table 2: Analysis of methods

| Parameters | Neural Network | Fuzzy Logic | SVM | Naïve Bayes | Decision Tree |
|---|---|---|---|---|---|
| Example Algorithms | Back propagation | Thresholds and weights applied on IF – THEN rules | Maximum & optimal margins by Gaussian theorem | Posterior Probability – Bayes Theorem | C4.5 , CART, J48 using splitting attribute entropy |
| Formula | Input Layer $w_{ij} = w_{ij} + \Delta w_{ij}$ Hidden Layer $w_{jk} = w_{jk} + \Delta w_{jk}$ | Fuzzy Set $\mu: X \to [0,1]$ | Margins Equations $w \cdot x - b = 1$ $w \cdot x - b = -1.$ | $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$ | Information Gain Gini Index |
| Advantages | Minimizes error in each level | Specification is obtained | Large data set is analyzed | Minimum error occurs | No domain knowledge is required |
| Disadvantage | Very slow working | Comparison increases | Range should be precise else outliers are observed | Multiple symptoms cannot handle and dependency in attributes | Selection of splitting attribute & over fitting |
| Approximate Accuracy | 80 - 90 % | 78 – 85 % | 85 – 90 % | 90 – 95 % | 94 – 96 % |

## V. PROPOSED ALGORITHM

### 5.1 Motivation

The proposed system is aims a enhancement in applications of tree distance function utilized. In other approaches similarity or distance score cannot be evaluated appropriately. The similarity function or score is evaluated with $S^*(T^P, T^q)$ .The evaluation is if $S^*(T^P, T^q)$ = 1 then trees are identical otherwise if $S^*(T^P, T^q) = 0$ then trees are totally distinct. If m is a constant number then similar nodes between $T^P$ and $T^q$ have strong similarity then they form an identical subtree mapping between $T^P$ and $T^q$. That is, identical subtree represents a similar sub structure between $T^P$ and $T^q$, where as m disjoint mapped nodes indicate no similar structure between two trees.

To conclude with in tree construction, the nodes are constructed on splitting attribute or the flag value as well as similarity score. Hence if continuous value is to be handled then it can prove to be fatal. Hence main motive behind this proposed system is to handle these problems.

### 5.2 Extended Sub Tree Similarity

The given $T^P and T^q$, proposed system of EST that is Extended Subtree, it helps to maintain tree structure by mapping subtrees of $T^P$ to the similar subtree of $T^q$. Now while mapping these $T^P and T^q$ there are some rules as follows

Rule 1: EST mapping means not only mapping only single nodes mapped together, but also identical subtrees mapped together.

Rule 2: No common subtrees $T^P$ and $T^q$ are allowed to mapped together, that is dissimilar trees cannot be mapped together. Rule 3 : In one to many mapping, a subtree of $T^P$ can be mapped into different subtrees of $T^q$ and vice versa.

Rule 4: m is the weighted as, $W(m_x) = (W(T^{px}) + WTqx/2$ where $WTpx$ and $WTqx$ are weights of subtrees in mapping. The $W(T^p)$ is calculated as

$$W(T^{px}) = \sum_{t_i^{px} \in T^{px}} W(t_i^{px}) \qquad \dots\dots(1)$$

where $W(T^{px})$ is the scalar unit, when $T^{px}$ is largest subtree that $t_i^{px}$ belongs to, and zero otherwise. Then we compute $S^*(T^P, T^q)$ based on all possible mappings such as

$$S(T^P, T^q) = \alpha \sqrt{\sum_{m_k \in M} \beta_k \, X \, W(m_k)^\alpha} \dots\dots (2)$$

where $\alpha, \alpha \geq 1$, is a coefficient to adjust the relation among different sizes of mappings. Then $\beta_k$ is the unit scalar, when the root nodes of $T^{Pk} and T^{qk}$ have same depth with respect to $T^P and T^q$ and it is equal to $\beta$ (a constant no between zero and one) otherwise, leads to enhancement of mapping of same depth regarding subtrees.

To normalize the similarity score, we divide it by its higher bound. Since $0 \leq \beta_k \leq 1$, we have $S(T^P, T^q) \leq \alpha \sqrt{\sum_{m_k \in M} W(m_k)^\alpha}$ . Further, $\alpha \sqrt{\sum_{m_k \in M} W(m_k)^\alpha} \leq \sum_{m_k \in M} W(m_k)$ where $\alpha \geq 1$ and $W(m_k)$ is a positive number. In addition, each node counted as one in weight calculation as, $\sum_{m_k \in M} W(m_k) \leq Max(|T^p|, |T^q|)$. This evaluates to, $S(T^P, T^q) \leq Max(|T^p|, |T^q|)$ and similarity function normalizes to,

$$S^*(T^P, T^q) = \frac{S(T^P, T^q)}{Max \ (|T^p|, |T^q|)} \dots\dots\dots\dots(3)$$

*5.3 Computational Algorithm*

Hypothesis $T_{i,j}^p$ represents a subtree of $T^p$ rooted to $t_i^p$ is mapped to identical subtree of $T^q$ rooted to $t_j^q$ namely $T_{j,i}^q$. Now evaluation of $S(T^P, T^q)$ is done in four steps as follows –

*Step 1: Identifying all mappings:* We evaluate all possible mappings, whether it may be valid or invalid (i.e invalid mappings will have weight zero from step 3 ), and store into two lists of nodes having each list for one each subtree. $T^p$ and $T^q$ are the inputs, while as $V^P$ and $V^q$ are the outputs (inputs for next step) . $V^P$ and $V^q$ are the two dimensional matrices where each element is a list of nodes represented as $V_{[i][j]}^p$ and $V_{[j][i]}^q$ to the list of nodes of mapped subtrees of $T_{i,j}^p$ and $T_{j,i}^q$ respectively.

In this step GetMapping(i,j) function results into two list of nodes ($V_{[i][j]}^p$ and $V_{[j][i]}^q$) for mapping. Its objective is to detect the largest mapping, which can be achieved at rooted children of $t_i^p$ and $t_j^q$. Now among these $t_i^p$ and $t_j^q$'s children, $t_{ia}^p$ is the $a^{th}$ child of $t_i^p$ node, where $1 \le a \le deg(t_i^p)$, and ia denotes index of $a^{th}$ child of $t_i^p$ node. Similarly $t_{jb}^q$ is the $b^{th}$ child of $t_j^q$ node, where $1 \le b \le deg(t_j^q)$, and jb denotes index of $b^{th}$ child of $t_j^q$ node. E is a matrix which indicates how children of $t_i^p$ and $t_j^q$ are matched. Also E is used to update $V_{[i][j]}^p$ and $V_{[j][i]}^q$. Therefore of $T_{i,j}^p$ and $T_{j,i}^q$ are identical so $\left|V_{[i][j]}^p\right| = \left|V_{[j][i]}^q\right|$.

*Step 2: Identifying each node's largest mapping:* A node $T^P$ or $T^q$ may belong to many mappings, so we consider largest sub tree in mapping for each node. To evaluate this, hypothesis of two arrays namely, $LS^p$ and $LS^q$ of size $T^P$ and $T^q$ respectively. $LS^p[i]$ indicates largest subtree that $t_i^p$ belongs to indexes of root nodes of mapping, denoted by $LS^p[i]_{mi}$ and $LS^p[i]_{mj}$ .The goal of this step is to fill $LS^p$ and $LS^q$ with appropriate values. Check if $\left|V_{[i][j]}^p\right|$ is larger than the subtree store it into $LS^p$ for that node and then update it as per the upliftment. Similarly follow for rach node in $V_{[j][i]}^q$ .

*Step 3: Compute the weight of each subtree:* For this step, evaluate $W(T_{i,j}^p)$ and $W(T_{j,i}^q)$ for all subtrees in mapping, which is stored into $W^p[i][j]$ and $W^q[j][i]$. If largest value as compared to previous value is found then add it to $LS^p$ and increment the weight of subtree. Similarly follow for $LS^q$ .

*Step 4: Calculate $S(T^P, T^q)$ :* In this step we have all subtree weights ($W^p$ and $W^q$) available. Then simply evaluate $S(T^P, T^q)$.



Figure 1. Pseudo code

## VI. RESULT ANALYSIS

This algorithm is implemented for diagnosing the heart diseases. The diagnosis conclude with the stage in which the disease is residing. The data is in continuous form, i.e range of values for every parameters is to be considered. There were 13 parameters to be considered for diagnosing the data. The Cleveland data set is been used for analysis purpose. The description of parameters can be given as follows:-

Table 3. Attributes for classification

| Sr No. | Parameter | Description |
|---|---|---|
| 1. | Age | |
| 2. | Gender | 0 – Female        1- Male |
| 3. | Chest Pain | 1- Typical Angina    2 – Atypical Angina<br>3- Non Angina Pain  4 – Asymptomatic |
| 4. | Treatbps | Resting Blood Pressure |
| 5. | Cholesterol | |
| 6. | FBS – Fasting Blood Sugar | 1- True          0 - False |
| 7. | RestECG – Resting Cardio graphic Results | 0 – normal<br>1 – having ST-T abnormality<br>2- probable or definite left ventricular hypertrophy. |
| 8. | Thalach – Maximum Heart Rate Achieved | |
| 9. | Exang – Exercise Induced Angina | 0 – Yes          1 - No |
| 10. | Oldpeak – exercise related to rest | |
| 11. | Slope – slope of peak exercise | 0 – up sloping    1- flat<br>2- down sloping |
| 12. | Cardiac arrest  - no.of major vessels by fluoroscopy | |
| 13. | Thal | 3- normal        6 – fixed Defect<br>7 – reversible defect |

By using these attributes the data is been classified to diagnose the disease of the patient. In ID3 algorithm, it cannot handle continuous data. Hence to implement ID3 algorithm, initially the data has to be converted into nominal form ie from continuous to non continuous form. After conversion it will evaluate the dataset to generate results. So for performing these steps, time required to evaluate comes around 54 ms. The time complexity required for evaluation comes around $O(n \log n)$. Also the accuracy achieved in this evaluation comes around 80.17 %. So to overcome the drawback of ID3 algorithm, extended sub tree approach was proposed. In this approach, continuous data can easily be handled and time required for evaluation is reduced. The time complexity of this algorithm is around $O(|T^p|, |T^q|) \times \min|T^p|, |T^q|$. Improvement in accuracy is observed as compared to ID3 algorithm.
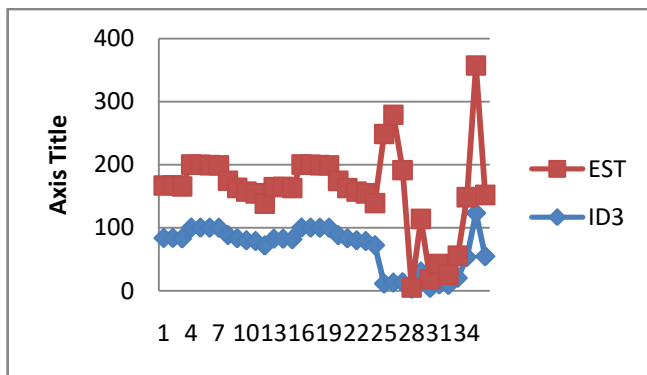


Figure 2. EST & ID3 Comparison

## VII. CONCLUSION

Clinical Decision Support System for heart diseases is very effective tool for diagnosing the diseases. It gathers the patient health information (PHI) and by using pre-determined algorithms or rules, it gives decisions. System will give decision of probability for patient been prone to heart diseases. Hence for implementation of such system Decision Tree technique will be an effective technique in classification. It is a simple tree like flowchart structure which helps in bifurcating the data in respective groups. The main goal of Decision Trees is in the intuitive representation that is easy to understand and comprehend. Also in decision tree construction, the nodes are constructed on splitting attribute or the flag value. Hence if continuous value is to be handled then it can prove to be fatal. Hence extension of sub tree is the approach to be implemented.

## ACKNOWLEDGMENT

## REFERENCES

[1]. Ali Shahbazi and James Miller"*Extended Subtree: A new similarity function for Tree structured Data*", IEEE Transactions on knowledge and Data Engineering, Vol 25, No 4, April 2014.

[2]. Christophe Damas, Bernard Lambeau, and Axel van Lamsweerde ,"*Analyzing Critical Decision-Based Processes* ", IEEE Transcations on Software Engineering, Vol. 40, No. 4 , April 2014

[3]. Tapuria,A., Austin, T., Sun, S., Lea, N., Iliffe S., Kalra, D., Ingram, D., Patterson, D., " *Clinical Advantages of decision support tool for anti coagulation control*", 2013 IEEE Point - of - Care Healthcare Technologies (PHT), Banglore, India, 16 - 18 January 2013.

[4]. Rodrigo C. Barros, M´arcio P. Basgalupp, Andr´e C. P. L. F. de Carvalho and Alex A. Freitas," *A Survey of Evolutionary Algorithms for Decision Tree Induction* ", IEEE Transcations on Systems, Man ,and Cybernetics - Part C: Applications & Reviews , Vol. X, No. X , January 2012

[5]. Gwenole Quellec, Mathieu Lamard, Lynda Bekri, Guy Cazuguel; "*Medical Case Retrieval From a Committee of Decision Trees*", IEEE Transactions On Information Technology in Biomedicine , Vol. 14, No. 5, September 2011

[6]. Kittipol Wisaeng ,"*Predict the Diagnosis of Heart Disease Using Feature Selection and k-Nearest Neighbor Algorithm* ",Applied Mathematical Sciences, Vol. 8, no. 83, 4103 - 4113, 2014.

[7]. Dhanashree S. Medhekar1, Mayur P. Bote2, Shruti D. Deshmukh, "*Heart Disease Prediction System using Naive Bayes*", INTERNATIONAL JOURNAL OF ENHANCED RESEARCH IN SCIENCE TECHNOLOGY & ENGINEERING VOL. 2 ISSUE 3, ISSN NO: 2319-7463, MARCH.-2013.

[8]. Shamsher Bahadur Patel 1, Pramod Kumar Yadav2, Dr. D. P.Shukla," *Predict the Diagnosis of Heart Disease Patients Using Classification Mining Techniques*", IOSR Journal of Agriculture and Veterinary Science (IOSR-JAVS) e-ISSN: 2319-2380, p-ISSN: 2319-2372. Volume 4, Issue 2, PP 61-64 , July - August 2013

[9]. Punam Suresh Pawar, D. R. Patil, "*Survey on clinical decision support system*", April 2012

[10]. Chaitrali S. Dangare Sulabha S. Apte, " *Improved Study of Heart Disease Prediction System using Data Mining Classification Techniques* ","International Journal of Computer Applications (0975 – 888) Volume 47– No.10, June 2012

[11]. R. S. Bichkar,Dipak V. Patil,"*Issues in Optimization of Decision Tree Learning*: A Survey", International Journal of Applied Information Systems (IJAIS) – ISSN : 2249-0868 Foundation of Computer Science FCS, New York, USA Volume 3 – No.5, July 2012

[12]. Hadi Sadoghi Yazdi,and Nima Salehi-Moghaddami," Multi *Branch Decision Tree: A New Splitting Criterion*",International Journal of Advanced Science and Technology Vol. 45, August, 2012

[13]. Miss. Chaitrali S. Dangare1, Dr. Mrs. Sulabha S. Apte, "*A DATA MINING APPROACH FOR PREDICTION OF HEART DISEASE USING NEURAL NETWORKS*", INTERNATIONAL JOURNAL OF COMPUTER ENGINEERING & TECHNOLOGY (IJCET), ISSN 0976 – 6375, Volume 3, Issue 3, pp. 30-40, October - December 2012

[14]. P.K. Anooj, "*Clinical decision support system: Risk level prediction of heart disease using weighted fuzzy rules*",Journal of King Saud University – Computer and Information Sciences , 2012

[15]. Yang Zhang,Simon Fong,Jinan Fiaidhi,and SabahMohammed, "*Real-Time Clinical Decision Support Systemwith Data*

*StreamMining*", Hindawi Publishing Corporation Journal of Biomedicine and Biotechnology Volume 2012

[16]. Mrs.G.Subbalakshmi,Mr. K. Ramesh, Mr. M. Chinna Rao ,"*Decision Support in Heart Disease Prediction System using Naive Bayes*" Indian Journal of Computer Science and Engineering (IJCSE), ISSN : 0976-5166 Vol. 2 No. 2 , Apr-May 2011

[17]. Sunita Pachekhiya, Gwalior,"*DISEASE DIAGNOSIS OF HEART MUSCLES USING ERROR BACK PROPAGATION NEURAL NETWORK*", International Journal of Engineering Science and Technology (IJEST), ISSN : 0975-5462 Vol. 3 No. 7, July 2011

[18]. Mai Shouman, Tim Turner, Rob Stocker,"*Using Decision Tree for Diagnosing Heart Disease Patients*", Proceedings of the 9-th Australasian Data Mining Conference (AusDM'11), Ballarat, Australia, CRPIT Volume 121 - Data Mining and Analytics 2011

[19]. Omar F. El-Gayar , Amit Deokar, Matthew Wills, "*Current Issues and Future Trends of Clinical Decision Support Systems (CDSS) *", IGI Global

[20]. Shawkat Ali , Kate A. Smith, "On learning algorithm selection for classification", Science Direct

[21]. Dawn DOWDING, Rebecca RANDELL, Natasha MITCHELL, Rebecca FOSTER, Carl THOMPSON, Valerie LATTIMER and Nicky CULLUM, "*Experience and Nurses Use of Computerised Decision Support Systems*", February 2010.

[22]. M. ANBARASI,E. ANUPRIYA, N.CH.S.N.IYENGAR, "Enhanced *Prediction of Heart Disease with Feature Subset Selection using Genetic Algorithm* ",International Journal of Engineering Science and Technology Vol. 2(10), 5370-5376, ISSN: 0975-5462, 2010

[23]. Shantakumar B.Patil, Y.S.Kumaraswamy, "*Intelligent and Effective Heart Attack Prediction System Using Data Mining and Artificial Neural Network* ", European Journal of Scientific Research ISSN 1450-216X Vol.31 No.4, pp.642-656, 2009

[24]. M.M.Abbasi, S. Kashiyarndi, "*Clinical Decision Support Systems: A discussion on different methodologies used in Health Care*"

[25]. Peter Szolovits,"*Uncertainty and Decisions in Medical Informatics*".

[26]. CARLA E., PAUL E. UTGOFF ,"*Multivariate Decision Trees*".

[27]. Omar F. El-Gayar, Amit Deokar, Matthew Wills," *Current Issues and Future Trends of Clinical Decision Support Systems (CDSS)*".