

## Machine Learning Techniques and Syntactic Pattern Recognition based Heart Disease Prediction for Smart Health

Prof. Shawni Dutta and Prof. Samir Kumar Bandyopadhyay  
Department of Computer Science, The Bhawanipur Education Society College, Kolkata, India.

### Abstract

Cardiovascular disease (CVD) may sometimes unexpected loss of life. It affects the heart and blood vessels of body. CVD plays an important factor of life since it may cause death of human. It is necessary to detect early of this disease for securing patients life. In this chapter two exclusively different methods are proposed for detection of heart disease. The first one is Pattern Recognition Approach with grammatical concept and the second one is machine learning approach.

In the syntactic pattern recognition approach initially ECG wave from different leads is decomposed into pattern primitive based on diagnostic criteria. These primitives are then used as terminals of the proposed grammar. Pattern primitives are then input to the grammar. The parsing table is created in a tabular form. It finally indicates the patient with any disease or normal. Here five diseases beside normal are considered.

Different Machine Learning (ML) approaches may be used for detecting patients with CVD and assisting health care systems also. These are useful for learning and utilizing the patterns discovered from large databases. It applies to a set of information in order to recognize underlying relationship patterns from the information set. It is basically a learning stage. Unknown incoming set of patterns can be tested using these methods. Due to its self-adaptive structure Deep Learning (DL) can process information with minimal processing time. DL exemplifies the use of neural network. A predictive model follows DL techniques for analyzing and assessing patients with heart disease. A hybrid approach based on Convolutional Layer and Gated-Recurrent Unit (GRU) are used in the paper for diagnosing the heart disease.

**Key Words:** Machine Learning, Deep Learning, Syntactic Pattern Recognition, Pattern Primitives and Heart Disease

### Introduction

The foremost reasons for high mortality rate over the globe are due to CVD. As per World Health Organization (WHO) statistics nearly 17.7 million people pass away every year in the globe [1-2]. Human heart along with blood vessels is known as Cardiovascular system [3].

Coronary artery disease (CAD), heart failure, cardiac arrest, and unexpected cardiac death are due to disorders of Cardiovascular system. It affects humans mostly due to uncontrolled behaviour in their daily life. The interior part of arteries of the heart consumes fatty deposits or plaque.). It is mainly cholesterol deposits within the arteries and it is known as atherosclerosis. These deposits may thicken and cause the coronary arteries to narrow. Due to this the amount of blood and oxygen flows in a reduced rate through the arteries to the heart. The narrowing of the arteries prevents blood and oxygen from flowing easily to the heart muscle. This effect will be happened to human as he/she grows ages.

Angina (pain, discomfort, or pressure in the chest) caused due to these symptoms. If blood flow is completely blocked by plaque or a blood clot that forms inside the narrowed coronary artery, a heart attack may occur. Coronary artery disease symptoms may include weight gain, weakness and fatigue, etc. of the patient [4].

Based on the above discussion, it can be inferred that CVD plays significant role in human's life. Early detection of this disease is necessary for saving patients life. CVD is often dependent on mental anxiety, daily lifestyle, working profile of people. Symptoms of anxiety, depression and stress may often lead to CVD [5]. For detection of CVD using two heterogeneous approaches such as syntactic pattern recognition based approach and predictive modeling using deep learning method.

Grammatical approach in the first process is used cardiac disease diagnosis [6]. In this method the patient data matrix was constructed initially [6]. It is used for classification of diseases. Based on diagnosis criteria, pattern primitives are identified. It is obtained from the updated diagnostic criteria published by American Heart Association and also from the medical literatures [7]. Based on patient data matrix, an input string is generated. One of the context free language i.e. Chomsky normal form are used to form the production rules for six diseases including Normal. For parsing the input string Cocke-Younger-Kasami (CYK) algorithm is used [6]. At the end, the parsing table will highlight the occurrence of disease.

In the second approach, an automated predictive model is favoured for CVD detection. Early heart disease can be predicted by utilising supervised machine learning approaches those takes patient's record as input. To explore the problem of heart disease detection, classification methods are implemented. It associates input variable for finding target classes based on training data. Attributes comprise of patient's details such as serum, cholesterol, etc. These features can form a good feature space while recognizing patients with cardiac symptoms. The proposed models acts for analyse the information of patients about their past health history records and predict their chances of affecting in cardiac trouble. This prediction will in turn benefit the doctors to provide well-versed decision and prescribe medicines and surgeries accordingly [8].

By means of machine learning approach, heart disease detection is focused in this chapter as one of the approaches. In order to diagnose CVD, it is necessary to extract knowledge from patient's health history database and identify relationship between interfering factors and heart disease probability. The proposed methods capture relevant health records of patient

and discovers the tendency of heart disease. Timely detection and screening play leading role in prevention of heart attacks. Deep learning (DL) [9] is implemented in this chapter for the heart trouble prediction by a means of medical data. Two models are exemplified for this purpose. This paper proposes Recurrent Neural Network(RNN)-based which assembles multiple Long Short Term Memory (LSTM) [10] layers where LSTM is known to be a variation of RNN. This neural network classifier receives all interfering factors as features and identifies patients with heart disease troubles. The second model consists of multiple GRU layers. For finding superior model a comparative study is drawn among the both specified models. Lastly the best model for CVD classification problem is selected on the comparative study.

## Related Works

CVDs are the principal reason of mortality worldwide per year that may reach an approximation of 23.6 million in 2030 [11]. The largest contributor of CVDs is Coronary heart disease (CHD). The damage of arterial wall is the main reason. The leading common indicator of CHD is Myocardial Infarction (MI). Angina pectoris is the former symptom of the pathology for 50% of patients [11]. Immediate diagnosis of CHD patients can save life. Image processing techniques can help early detection of CHD.

For heart disease detection initially the electrocardiogram (ECG) is performed. It was started late in the 1950's. The diagnosis of disease is made by researchers using non-syntactic methods as well as syntactic methods and hybrid methods [6, 12]. The syntactic method is used for analyzing ECG pattern. This method is not much used in pattern analysis and a few works have been done till date. Only specific aspects of these areas are looked upon by researchers. For peak recognition in ECG's using Context-free grammar is described in [12].

A pattern in syntactic approach is considered to have a complex construction, which is decomposed into sub-patterns that in turn are decomposed into simpler sub-patterns, etc. In cardiology an ECG signal pattern is also treated as a linear structure, which consists of separable substructures describing the different phases of human heart's beating (e.g. P wave, T wave, ST segment, QRS complex), A set of various structures is perceived as a formal language. Words (structural patterns) are analyzed by formal automata which not only are able to identify proper categories (diseases) for patterns, but also can characterize their structural features. Therefore, syntactic pattern recognition seems to be convenient, if a descriptive structural characterization is a goal of ECG analysis rather than only its classification (i.e. assigning an ECG signal to one of classes of heart dysfunction phenomena) [13]. Electrocardiogram (ECG) is often utilized as common but vital sign from the clinical environment perspective. Analyzing ECG often reveals many cardiac disorders.

Existing literatures on automatic ECG classification are clubbed into different clusters for the review of classification process. In ECG-based computer-aided-diagnosis system unwanted information in ECG waves, parts of ECG wave detection, heartbeat classification, etc, are necessarily removed for proper diagnosis of disease. Here two approaches are discussed. One

is grammar based classification of diseases and other is machine learning based hybrid approach for classification of diseases [14-15].

Machine Learning (ML), specialized field of AI, can be used in healthcare that analyzes numerous different data points, recommends outcomes, provides well-timed risk scores, defined resource allocation, and delivers many other applications. The opportunities for improving clinical decision support can be made by ML. ML techniques are often related data mining procedure. From the data mining point of view, ML techniques can be said that data mining examines an enormous amount of data and sets a particular outcome based on those examined data. ML focuses on achieving that goal by using harvested data for modeling smart intelligent automated tool. By implementing data mining rules, data related to coronary illness is extracted from a large database. For this purpose, weighted association implemented in [16]. Using rule mining algorithms on patients' dataset, heart disease is predicted. Prediction results achieved 61% training accuracy and 53% testing accuracy.

Historical medical data is utilised in order to predict Heart Disease using ML techniques [17]. 462 instances of South African Heart Disease dataset used for prediction purpose. All these algorithms used validation method. It is 10-fold cross validation. The probabilistic Naïve bayes classifier performed better in comparison to other classifiers [17].

Heart Failure (HF) is classified into categories such as HF with preserved ejection fraction (HFPEF) and HF with reduced ejection fraction (HFREF) [18]. Various classification methods are used for detecting patients with heart failure. Several classification methods such as classification trees, random forests, bagged classification trees, boosted classification trees, and SVMs and for prediction, logistic regression, regression trees, bagged regression trees, random forests, and boosted regression trees are utilised for detecting patients with aforementioned three categories of heart failure. These are tree-based methods and regression trees for predicting and classifying HF subtypes.

K. Gomathi et al [19] predicted heart disease using Naïve Bayes Classifier and J48 classifier. They have concluded that Naïve Bayes classifier reaches an accuracy of 79% where J48 classifier reaches an accuracy of 77%. P.Sai Chandrasekhar Reddy et al used ANN for predicting Heart disease by considering relevant features such as heart rate, blood pressure etc [20]. Boshra Brahmi et. al. [21] employed several classification techniques such as J48, KNN, SMO, and Naïve Bayes for diagnosing heart disease. Instead of focusing on feature selection, emphasis is given on all relevant features for heart disease diagnosis and prediction [22]. This prediction modeling is implemented by assembling Random Forest with a linear model. Another study considered Arrhythmia which is irregular changes of normal heart rhythm as a prediction field [23]. Arrhythmia prediction is accompanied by implementing CNN which accepts ECG signals as input.

## Datasets

Datasets having ECG waves are collected from hospitals of West Bengal. The middle aged people with the range from 40 to 70 are considered. Others are taken from the American Heart Association [12].

This study implements deep learning based study for implementing computer aided classification. UCI machine learning repository is used for predicting cardiac disorder of a patient. Various attributes are in the dataset [24]. However, the attribute ‘target’ is utilized as output class of the prediction. Figure1 presents the overall histogram representation of the dataset. For obtaining a balanced dataset, preprocessing techniques are performed. After collecting the dataset some pre-processing techniques such as NaN values handling, scaling and transformation of some attributes such as age, cholesterol level etc are performed. This will assist the classifier in obtaining better predictive results. This pre-processed data is divided into 67:33 as training and testing dataset. Training data is given as input to the classifier model for learning process and after that testing dataset is used for obtaining prediction results. The distribution of cardiac and non-cardiac patients on the dataset is shown in Figure 2 and Table 1.

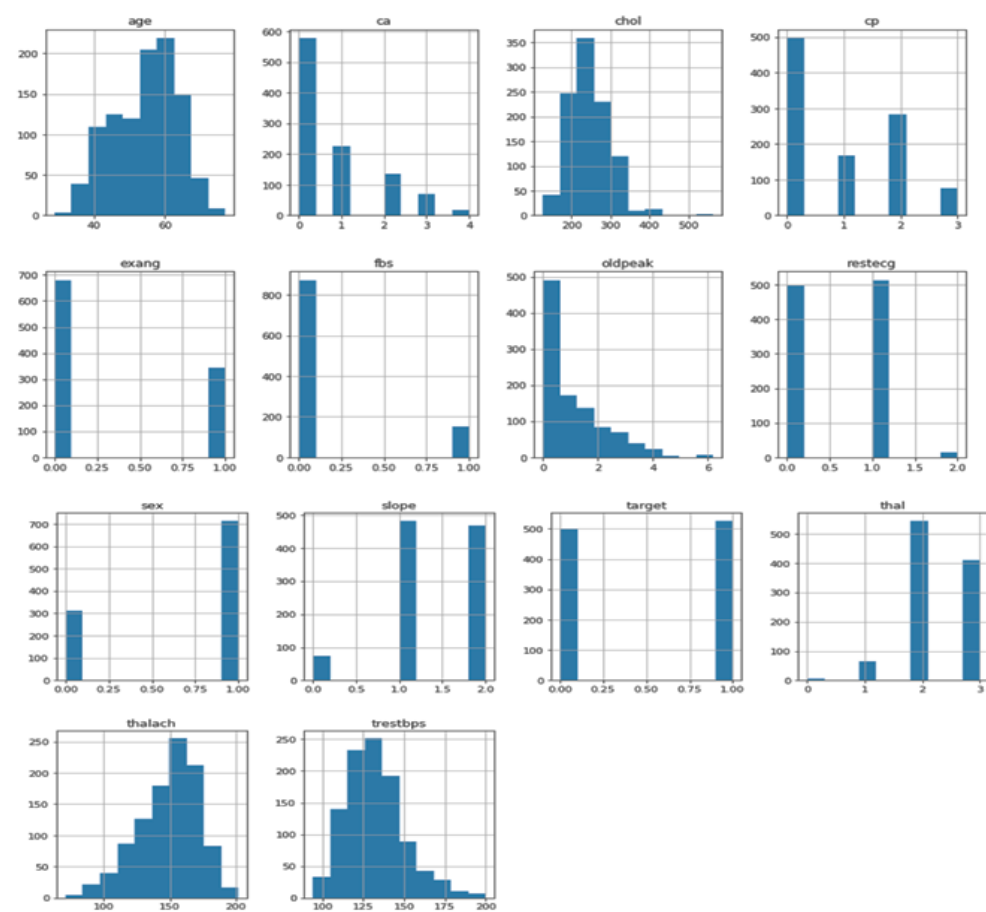


Figure 1: Histogram interpretation of Cardiac disease dataset

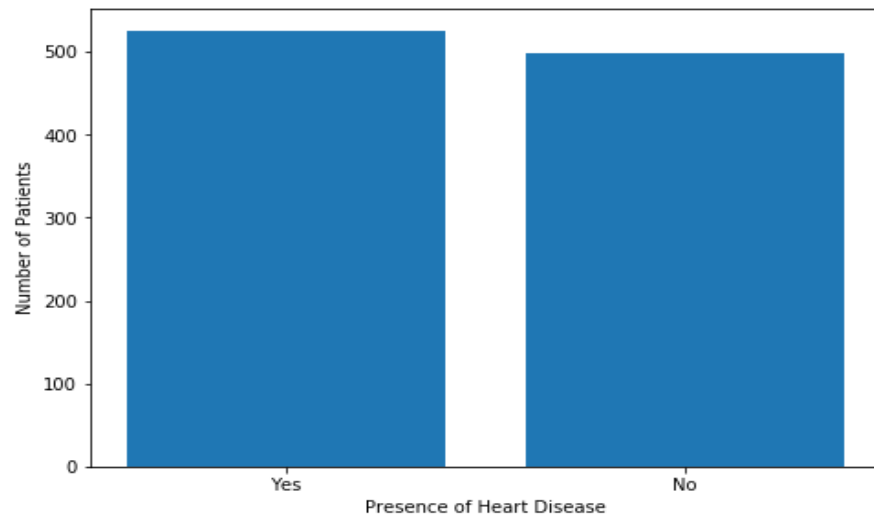


Figure 2: Distribution of target attribute over the collected dataset.

Table 1: Understanding of the Heart disease dataset

Attribute (Explanation)	Attribute Type	Values
Age	Numeric	40-70
Sex	Categorical	0-female, 1-male
Cp (Chest Pain)	Categorical	0: asymptomatic, 1: atypical angina, 2: non-anginal pain, 3: typical angina
Trestbps (The patients' resting blood pressure of the patient during the admission time; measured in the unit of mm Hg)	Numeric	94-200
Chol: Measurement of cholesterol in mg/dl	Numeric	126-564
fbs: fasting blood sugar of the patient measured in mg/dl	Binary	(if fbs > 120 mg/dl, 1 = true; otherwise 0 = false)
restecg: indicates the resting electrocardiographic outcomes	Categorical	0: indicates possibility of left ventricular hypertrophy using Estes' criteria 1: normal Two having ST-T wave abnormality (T wave inversions and/or ST elevation or depression of > 0.05 mV)
thalach: highest heart rate	Numeric	71-202

observed		
exang: Exercise induced angina	Categorical	0: no 1: yes
Oldpeak: ST depression induced by exercise relative to rest	Numeric	0.0-6.2
ST segment slope	Categorical	It varies from 0 to 2 depending on down, flat and up sloping
ca: quantity of major vessels	Categorical	1-4
Thalassemia (thal)	Categorical	Zero: indicates NULL (no Thalassemia) One: indicates fixed defect as some portion Two: Normal blood flow is denoted by two Three: Abnormal blood flow is observed
Presence of Heart disease (target)	Binary	1/0 or yes/no

## Proposed Methods

### Proposed Method 1: Pattern Recognition with Syntactic Recognition based Approach

While detecting coronary artery disease, contaminated recordings create major problem. So, pre-processing steps are highly recommended before processing ECG waves. For removing noise low-pass filter as well as high-pass filter is used. For power source interference 50 Hz notch filter is used. Figure 3 shows a normal ECG wave. The ECG wave is taken through 12 lead systems. Six electrodes are placed on the limbs. On chest six are placed.

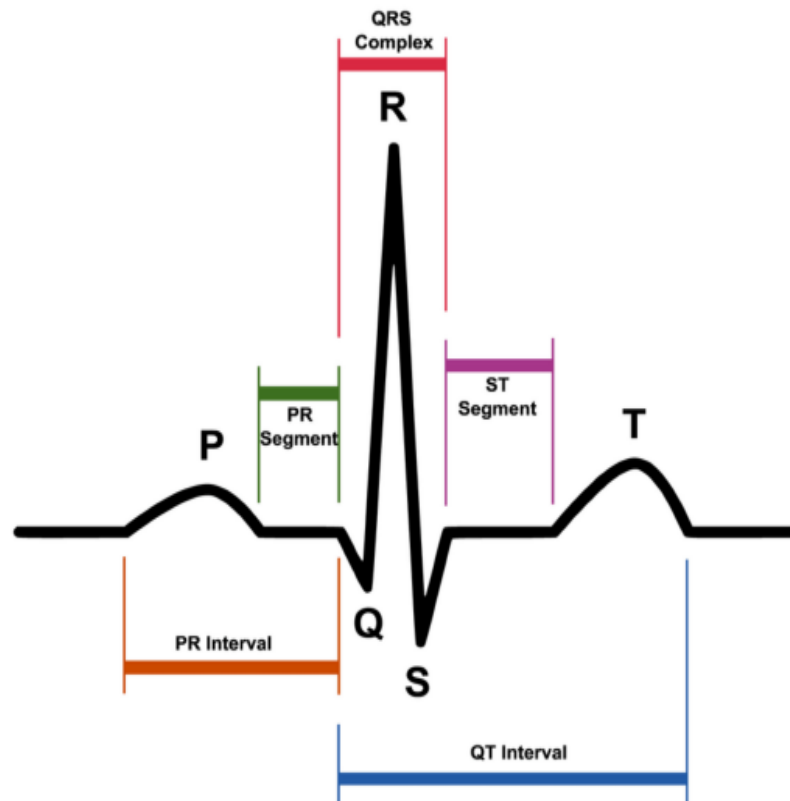


Figure 3: Normal ECG waveform and its feature patterns [25].

P, Q, R, S, and T waves are called PQRST complex. “R-R interval”, corresponds to a cardiac cycle. The following parameters are also used for diagnosis of heart disease from ECG Wave.

1. The end of P wave to the beginning of Q wave is denoted by PR interval. The starting of P to the end of Q is specified as PQ interval. The horizontal portion from the end of P wave to the beginning of Q wave is known as PR segment. The depolarization wave is identified by this segment.
2. The end of S wave to the start of T wave is ST segment.
3. The starting of Q wave to the end of T wave is called QT interval [25].

The syntactic methods of pattern recognition for cardiac diseases diagnosis is the main aim of the first approach [6]. Initially for generating the grammar in CYK normal form patient data matrix is used. The grammar is context free since there is no dependency of consecutive pattern primitives. These primitives are treated as terminals of the grammar. Non-terminals are created based on the terminals for classification of heart diseases. Here pattern primitives are terminals of the grammar. The updated diagnostic criteria published by the American



Heart Association and from a review of the medical literatures are used as diagnostic criteria.[14-15]. Initial string is based on the patterns primitives and patient data matrix. It is always required that the input string is parsed by the Production Rules of the grammar. Production rules of grammar are generated using terminals and non-terminals of the grammar. Production rules are in the form of CYK normal form. It is as per definition of context free language and it is developed by using diagnostic rules for describing five cardiac diseases besides the normal ECG. The Cocke-Younger-Kasami algorithm is used to parse the input string [6]. If the patient ECG has a sign of any abnormalities then the first column of the top of the parsing table will show the disease. Right- and left-bundle branch block, left-ventricular hypertrophy, left-anterior hemi block, and left-atrial hypertrophy-hereafter abbreviated to RBBB, LBBB, LVH, LAN and LAHI are chosen as the five diseases. Normal ECG wave is considered as six one for classification of disease. Four different areas of the ventricles are taken since these represents one of the vital part of the heart. Left-atrial hypertrophy is considered since it is common. It is aimed here for diagnosing a patient with normal symptom or abnormal symptom. For simplicity the five diseases are denoted as D1, D2, D3, D4 and D5. The selection of primitive selection is both problem-oriented and pattern-dependent. There is no general solution to this problem as yet [12]. The diseases and abnormal findings are identified as relationship for forming decisions regarding pattern and it is then transformed into pattern primitives. The patient data matrix is checked against Diagnostic criteria and patient data matrix are checked and transformed it into binary decision i.e. satisfied primitive and unsatisfied primitives. Either of two types is based on whether the particular condition is satisfied or not. A set of ten primitives in terms of notations A1, A2, B, C and D have been selected for the five diseases.

### D1:

- (A1) The Amplitude of R wave ( $A_R$ ), amplitude of S wave ( $A_S$ ) and amplitude of R' wave, of each complex is greater than 6 mm in lead  $V_1$  or  $V_2$ . Also the width of R' ( $D_R$ ) is greater than 0.025 s in lead  $V_1$  or  $V_2$ .
- (A2) Ventricular activation time ( $T_{VA}$ ) is greater than 0.44s in  $V_1$  or  $V_2$ .
- (B) Duration of S wave ( $D_S$ ) in lead I is greater than or equal to 0.03 s.
- (C) QRS duration ( $T_{QRS}$ ) is greater than 0.12 s.

### D2:

- (A) QRS duration ( $T_{QRS}$ ) is greater than 0.12s.
- (B)  $A_R$ ,  $A_S$  and  $A_R$  each complex is greater than 6mm in at least one of the leads I, AVL,  $V_5$  and  $V_6$  or notched R wave (i.e. the duration of R wave is greater than 0.44s) present in at least one of the leads I, AVL,  $V_5$  or  $V_6$ .

### D3:

- (A 1)  $A_R$  is greater than 27 mm in lead  $V_5$  or  $V_6$ .
- (A 2) Q wave amplitude  $A_Q$  or  $A_S$  in lead  $V_1$  plus  $A_R$  in lead  $V_5$  or  $V_6$  is greater than or equal to 35 mm.

- (A3)  $A_R$  is greater than or equal to 13 mm in lead AVL.  
 (A4)  $A_R$  in lead I plus  $A_S$  in lead II is greater than or equal to 26mm.  
 (B) Patient's age is above 30 years.

**D4:**

- (A) Left-axis deviation (LAD) is between  $-45^\circ$  and  $-60^\circ$ .  
 (B) Q wave duration is less than or equal to 0.02s in lead I and aVL.  
 (C)  $A_{VF}$  is less than 5 mm in leads I, II, III and AVF.  
 (D) Normal QRS duration (pure LAH can increase the QRS duration no more than 0.02 s, thus a QRS duration of 0.1 s indicates the coexistence of RBBB or some other form of ventricular conduction abnormality).

**D5:**

- (A 1) Notched P wave amplitude ( $A_P$ ) and P' wave amplitude ( $A_P$ ) is greater than 1 mm in leads I, II or AVL.  
 (A 2)  $A_P$  is greater than 3 mm in lead I or in lead AVL, or equal to 3.5 mm in lead II.  
 (B) Overall P wave duration ( $D_P$ ) is greater than 0.11 s.  
 Considering the diagnostic criteria as specified above the pattern primitives are selected based on whether the diagnostic criteria is satisfied or not. This is shown in Table 2.

Table 2: Diagnostic criteria and corresponding primitive notation suggested

Disease Name	Diagnostic Criteria	Notation used for 'satisfied' primitives	Notation used for 'unsatisfied' primitives
(1)	(2)	(3)	(4)
RBBB	(A 1)	B	M
	(A 2)	b	m
	(B)	b	m
	(C)	b	m
LBBB	(A)	c	k
	(B)	c	k
LVH	(A 1)	e	l
	(A 2)	e	l
	(A 3)	e	l
	(A 4)	e	l
	(B)	e	l
LAH	(A)	h	j
	(B)	h	j
	(C)	h	j
	(D)	h	j
LAH1	(A 1)	A	N
	(A 2)	a	n
	(B)	a	n

An input string is generated based on the diagnostic pattern primitives. This input string defines complete characterization of the disease structure taken into consideration. This representation will produce a disease pattern that comprises of basic elements those are

related the disease present in the input ECG wave. The string generation algorithm is described below.

### Algorithm

**Input:** Patient data matrix.

**Output:** A string Z of symbols is formed by the alphabet set = {a, b, c, e, h, j, k, l, m, n}. The symbols indicate pattern primitives either 'satisfied' or 'unsatisfied' condition for the considered disease.

Step 1. Assign  $p = 1$ .

Step 2. Assign  $q = 1$ .

Step 3. If the condition 'q' for the disease p is satisfied, then consider  $z_p =$  satisfied primitive and go to Step 5.

Step 4. Set  $z_p =$  unsatisfied primitive.

Step 5. If q is greater than the total number of criteria in disease p, then go to Step 7.

Step 6. Set  $q = q + 1$  and go to Step 3.

Step 7. Now concatenate  $z_p$  with  $z_{p-1}$ , to form the complete string z.

Step 8. If p is greater than the total number of considered diseases, then continue, otherwise set  $p = p + 1$  and go to Step 2.

Step 9. Exit.

Assume that the sample electrocardiogram is retrieved from a 58-year-old male patient. It has already been processed by using the described algorithm and this will yield the patient data matrix which is presented in Table 3. To shorten the discussion, we denote the following partial information from Table 3.

$$(1) \ A_{Rv1}, A_{Sv1}, A_{R'v1}/A_{Rv2}, A_{Sv2}, A_{R'v2} < 6 \text{ mm.} \\ D_{R'v1}/D_{R'v2} < 0.025 \text{ s.}$$

$$(2) \ T_{VAv1} < 0.04 \text{ s.}$$

$$(3) \ D_S < 0.03 \text{ s.}$$

$$(4) \ T_{QRS} < 0.12 \text{ s.}$$

⋮

After collecting the above information, we obtained the following string using Table 2.

m m m m.....

Table 3: Patient Data Matrix

Parameters	I	II	III	$A_{VR}$	$A_{VL}$	$A_{VF}$	$V_1$	$V_2$	$V_3$	$V_4$	$V_5$	$V_6$
PA	2.00	2.00	-2.00	-1.00	1.00	1.00	-1.00	2.00	2.00	1.00	1.00	1.00
PD	0.12	0.12	0.11	0.12	0.12	0.12	0.12	0.12	0.12	0.12	0.12	0.12
P'A	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
P'D	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
QA	-1.00	-1.00	0.00	0.00	0.00	-1.00	0.00	0.00	0.00	-2.00	-2.00	-2.00
QD	0.02	0.02	0.00	0.00	0.00	0.02	0.00	0.00	0.00	0.02	0.02	0.02
RA	20.00	11.00	0.00	0.00	17.00	3.00	0.00	3.00	3.00	40.00	40.00	40.00
RD	0.08	0.09	0.00	0.00	0.01	0.05	0.00	0.02	0.02	0.08	0.08	0.08
R'A	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
R'D	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
SA	-2.00	0.00	-15.00	-18.00	0.00	-11.00	-37.00	-22.00	-22.00	-12.00	-12.00	-12.00
SD	0.02	0.00	0.11	0.12	0.00	0.06	0.12	0.10	0.10	0.02	0.02	0.02
S'A	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
S'D	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
TA	2.00	3.00	1.00	-3.00	-3.00	2.00	-7.00	6.00	6.00	8.00	8.00	8.00
TD	0.12	0.12	0.11	0.12	0.12	0.12	0.12	0.12	0.12	0.12	0.12	0.12
VAT	0.06	0.06	0.00	0.00	0.07	0.04	0.00	0.01	0.01	0.07	0.07	0.07
PR	0.16	0.16	0.14	0.16	0.16	0.16	0.16	0.16	0.16	0.16	0.16	0.16
QT	0.28	0.28	0.25	0.28	0.28	0.28	0.28	0.28	0.28	0.28	0.28	0.28
ST	0.05	0.06	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
QRS	0.11	0.10	0.11	0.12	0.11	0.12	0.12	0.12	0.12	0.12	0.12	0.12
S-TON	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	58 AGE	-21.00 AXIS	5 HEART RATE									

After the string generation operation is completed, the immediate task to be accomplished is to specify diseases by a means of syntax analysis. The proficiency of a syntax analyser is dependent mainly on the grammar that generates the language and also depends on the parser that evaluates the syntactic correctness of an input string. In order to describe the considered normal and disease patterns, Context-free language in Chomsky normal form has been utilized. The names of the diseases with other symbols are taken as non-terminals. The names used for such non-terminals are the same as those used in conventional ECG nomenclature so that they can be easily understood. The diagnosis grammar describing the normal as well as the five disease patterns is given below.

$$A_4 = L \ E/E \ L/A_3 \ E/A_4 \ L/E \ E/A_5 \ L$$

$$L = l$$

$$E = e$$

$$A_3 = A_3 \ L/L \ L$$

$$A_5 = A_4 \ E/A_5 \ E$$

$$J = j$$

$$H = h$$

$$A_7 = H \ H$$

$$A_8 = A_7 \ H$$

$$A_9 = A_8 \ H$$

$$LBBB = C \ C/RBBB \ LBBB/NORMAL \ LBBB/LBBB \ A_4/LBBB \ A_3/LBBB \\ A_5/LBBB \ E/LBBB \ H$$

$$LVH = NORMAL \ A_5/NORMAL \ LVH$$

$$LAH = NORMAL \ A_9/LVH \ A_9/NORMAL \ H/LAH \ H/NORMAL \ LAH/RBBB \\ LAH/LVH \ H$$

$$LAH1 = F \ A$$

$$F = NA/AN/AA$$

$$A = a$$

$$N = n$$

The production rules are in the meta language BNF (Backus normal form). Many variants of BNF are in use. WSN (Wirth syntax notation) is used for convenience [32]. It is important to know whether the string belongs to  $L_{\text{Diagnosis}}^G$  or not.

The proposed method uses the Cocke-Younger-Kasami (C-Y-K) bottom-up parsing algorithm and it produces a structured table. The tabular form is well suited for Physicians to have a quick look about the condition of patient. The patient data matrix is now validated against diagnostic criterias and a string comprising of primitives is formed. Suppose the input string is composed of the following:

$$x = m^4 k c e^3 I e j^2 n^2 a$$

'm' means single occurrence and 'm<sup>2</sup>' represents the occurrence of primitive m twice and so on. Table 4 shows the parsing table. It indicates the disease or normal as per convention used in the conventional parsing table. The occurrence of diseases is investigated on inspection of the first column and particular location namely the third row, (total number of diagnostic criteria present in N-number of diseases-the number of criteria present in LAHI) plus one of the Parsing Table. 'NORMAL' indicates the left atrial hypertrophy diseases are not present in

Table 4: Parsing Table

NORMAL																			
NORMAL	NORMAL																		
NORMAL	NORMAL	NORMAL																	
NORMAL	NORMAL	NORMAL	0	NORMAL															
NORMAL	NORMAL	NORMAL	0	NORMAL															
NORMAL	NORMAL	NORMAL	0	NORMAL															
LVH	NORMAL	NORMAL	0	NORMAL															
NORMAL	LVH	NORMAL	0	NORMAL															
LVH	NORMAL	LVH	0	NORMAL															
NORMAL	LVH	NORMAL	0	NORMAL															
0	NORMAL	LVH	0	LVH	0	0	0	0	0	0	NORMAL								
NORMAL	0	NORMAL	0	NORMAL							NORMAL	NORMAL							
0	NORMAL	0	0	LVH	A <sub>3</sub>						NORMAL	NORMAL	NORMAL						
NORMAL	0	NORMAL	0	NORMAL	A <sub>4</sub>	A <sub>5</sub>					NORMAL	NORMAL	NORMAL						
NORMAL	NORMAL	0	0	0	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>				NORMAL	NORMAL	NORMAL	NORMAL					
NORMAL	NORMAL	NORMAL	0	NORMAL	A <sub>4</sub>	A <sub>4</sub>	A <sub>4</sub>	A <sub>4</sub>			NORMAL	NORMAL	NORMAL	NORMAL	NORMAL	F			
M	M	M	M	K	C	E	E	E	L	E	J	H	J	J	N	N			
m	m	m	m	k	c	e	e	e	l	e	j	h	J	j	n	n			

Deep Learning (DL) is a specialized area of Machine Learning (ML) which enforces automatic learning of abstract information from large database without incorporating manual feature engineering methods. Deep neural networks are capable to compute complex functions by extracting features from input data. These computations are dependent on number of hidden layers and other parameters. For accompanying the complex computations, activation functions are used. Activation functions are advantageous in executing complicated computations and associates input signal into output signal within a certain range [26, 33-38].

Over-fitting is a serious problem that is faced mostly by neural network based model. This problem occurs when a model learns noise present in the training data which in turn negatively impacts the efficiency of the model on unknown data. This problem can be

eliminated by incorporating drop out layers. During each of the training iterations, Dropout layer randomly deactivates a fraction of the units or connections in a network [27]. Once the neural model is configured, it undergoes through a training process. The training process is executed through one cycle which is known as an epoch. In an epoch, the dataset is partitioned into smaller sections. For completing execution of each epoch, an iterative process is carried out by a means of batch size that considers subsections of training dataset for completing epoch execution [28]. The training process is also accompanied by a training criterion, known as binary cross entropy function as a binary classification problem is implemented in this study. Binary cross entropy finds out the difference between the true value (which is either 0 or 1) and the prediction for each of the classes and then class-errors are averaged out to measure the final loss [29].

Any machine learning models depend on some predefined metrics such as accuracy, precision, recall, f1-score, MSE and cohen-kappa statistics. These metrics help in identifying the best problem solving tactic. Accuracy [30] determines the percentage of true predictions over the whole number of instances considered. However, accuracy evaluation may not be enough since it does not reflect wrong predicted cases. For resolving the above mentioned problem, two more metrics known as, Recall and Precision can be yielded. Precision [30] ascertains the fraction of correct positive results over the number of positive results predicted by the classifier. The number of correct positive results divided by the number of all relevant samples is measured by recall [30]. F1-Score or F-measure [30] is another parameter that is basically the harmonic mean of both precision and recall. Mean Squared Error (MSE) [30] is another evaluating metric that can differentiate the prediction observation from actual observation of the test samples. A model having higher values of accuracy, F1-Score and lower MSE value indicate best problem-solving technique. Cohen-Kappa Score [31] is a statistical parameter that discovers inter-rater agreement for qualitative items for classification technique.

The objective of any classifier model is to map input variables into target variables considering the training dataset. The proposed classifier employs deep learning techniques in order to recognize whether a patient has heart disease or not. The proposed method uses LSTM-BRNN model for such prediction. A stacked LSTM-BRNN model is implemented as the second approach that stacks four Bidirectional LSTM layers and four dense layers. This stacked LSTM-BRNN model is built up using 256, 128, 64, 16 nodes respectively in every single layer. To avoid over-fitting problem, each layer is incorporated with 20% of dropout regularization. Next, four fully connected layers are stacked by including 8,4,2,1 number of nodes respectively. The first four LSTM layers and the final dense layers are activated using sigmoid activation function. Finally, the above mentioned layers are assembled using “adam” optimizer. This model is accompanied by binary cross entropy loss function. Construction of this model is dependent on epoch size of 100 and batch size of 32. The mentioned hyper-parameters have undergone through a series of possible values and the mentioned values are picked up. This fine-tuning operation will support in attaining the best problem-solving approach. Once this model is constructed, training data is fitted into the proposed model.



During the training phase, the presented neural network model accepts a total of trainable 1,367,993 parameters to retrieve prediction. An in depth description in terms of layers, type of layers, activation function used, output shape produced by each layer, number of parameter accepted by each layer is summarised in Table 5. The proposed model consists total of 12 layers, out of which 4 layers are of LSTM neural network.

The same configuration is used by stacked bi-directional GRU model. Table 6 describes the detailed construction of the second model. Description of all the hyper-parameters for Stacked Bidirectional LSTM model as well as Stacked Bidirectional GRU model is summarized in Table 5 and Table 6 respectively.

Table 5: Stacked bidirectional LSTM model's description

Layer	Number of Nodes/Percentage Rate	Output Shape	Number of Parameters Received	Activation function Used
Bidirectional LSTM layer	256	(None, 30, 512)	528384	Sigmoid
Dropout Layer	20%	(None, 30, 512)	0	None
Bidirectional LSTM layer	128	(None, 30, 256)	656384	Sigmoid
Dropout Layer	20%	(None, 30, 256)	0	None
Bidirectional LSTM layer	64	(None, 30, 128)	164352	Sigmoid
Dropout Layer	20%	(None, 30, 128)	0	None
Bidirectional LSTM layer	16	(None, 32)	18560	Sigmoid
Dropout Layer	20%	(None, 32)	0	None
Dense layer	8	(None, 8)	264	None
Dense layer	4	(None, 4)	36	None
Dense layer	2	(None, 2)	10	None
Dense layer	1	(None, 1)	1	Sigmoid

Table 6: Stacked bidirectional GRU model's description

Layer	Number of Nodes/Percentage Rate	Output Shape Obtained from each layer	Number of Parameters Received	Activation function Used
Bidirectional GRU layer	256	(None, 30, 512)	508384	Sigmoid
Dropout Layer	20%	(None, 30, 512)	0	None
Bidirectional GRU layer	128	(None, 30, 256)	553394	Sigmoid
Dropout Layer	20%	(None, 30, 256)	0	None
Bidirectional GRU layer	64	(None, 30, 128)	124852	Sigmoid
Dropout Layer	20%	(None, 30, 128)	0	None
Bidirectional GRU layer	16	(None, 32)	17690	Sigmoid
Dropout Layer	20%	(None, 32)	0	None
Dense layer	8	(None, 8)	249	None
Dense layer	4	(None, 4)	34	None
Dense layer	2	(None, 2)	9	None
Dense layer	1	(None, 1)	1	Sigmoid



Table 7: Best Hyper-parameter Specification for Stacked bidirectional LSTM model

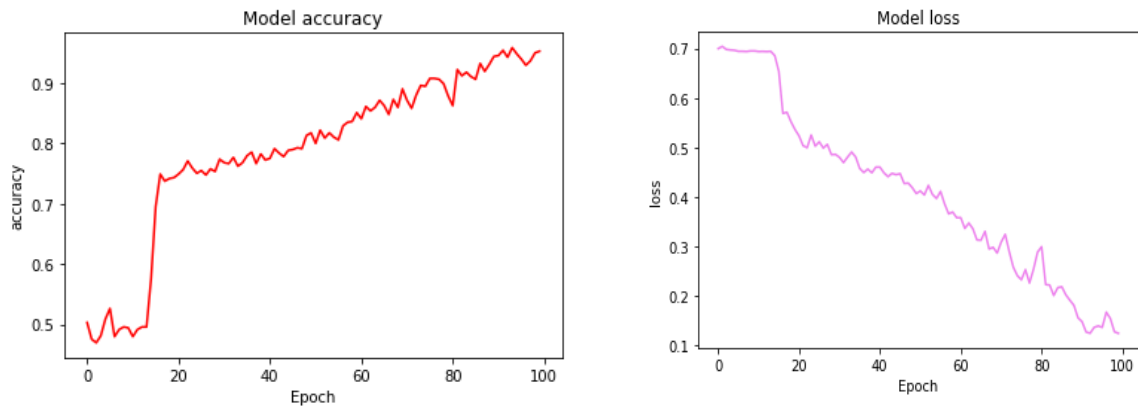
Hyper-parameters Used	Values
Number of Epochs	100
Optimizer Used	Adam
Loss Function	Cross-Entropy
Batch Size	64

Table 8: Best Hyper-parameter Specification for Stacked bidirectional GRU model

Hyper-parameters Used	Values
Number of Epochs	50
Optimizer Used	Adam
Loss Function	Cross-Entropy
Batch Size	32

During training process of the Stacked Bidirectional LSTM model, accuracy and loss are calculated for each epoch as depicted in Figure 4. As the quantity of epochs grows, the accuracy increases gradually and reaches around a value of 0.95. In contrast, the loss gradually decreases and attains lowest value around 0.12. Once the training process is done i.e., after completing 100 epochs, accuracy, f1-score, cohen-kappa score and MSE rate for unlabelled dataset. Table 7 provides the prediction efficiency for the presented model. It is to be noted that the proposed stacked bi-directional LSTM model has 4 LSTM layers. In Table 8 it is also shown as the model efficiency is increased over 1, 2 and 3 LSTM layers. Increasing more than 4 LSTM layers is not enhancing much substantial efficiency. Hence it is restricted to have 4 LSTM layers as model component.

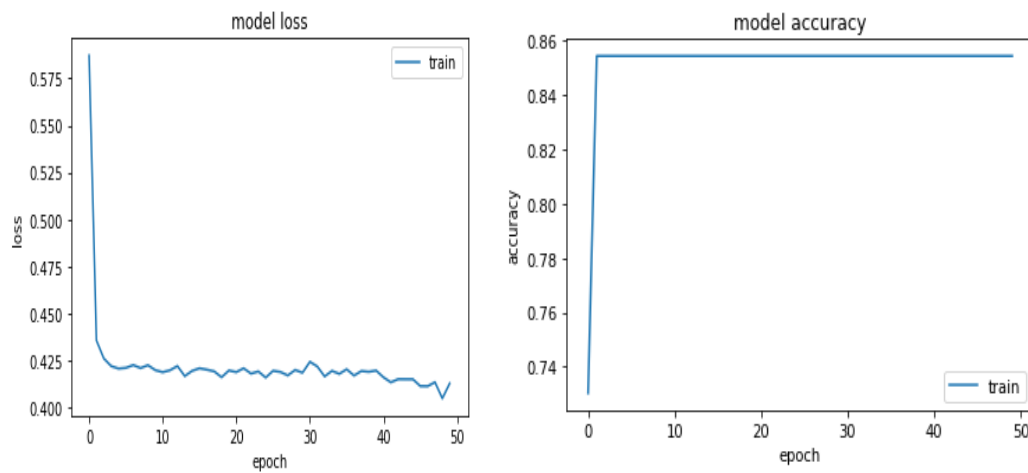
As shown in Figure 5, the Stacked Bidirectional GRU model is trained for 50 epochs. Increasing the number of epochs more than 50 is not contributing the efficiency of the model. Hence, it is restricted till 50 epoch size. The training loss declines rapidly within 10 epochs and later decreases gradually as number of epochs increases. After 50 th epoch it approaches a loss of 0.394. During training this model starts from obtaining a lower value accuracy which is increased till 0.8542 after certain epochs. Table 9 provides the performance of prediction for the proposed GRU based model. Comparative study among the 1, 2, 3 and 4 GRU layers is also described in Table 10.



(a) Model Accuracy.

(b) Model loss

Figure 4: Training process of Stacked Bi-directional LSTM model.



(a) Model Accuracy.

(b) Model loss

Figure 5: Training process of Stacked Bi-directional GRU model

Table 9: Performance of prediction drawn by StackedBi-directionalLSTM model.

Number of Layer Used	Accuracy	Cohen-Kappa Score	F1-Score	MSE
4	93.22%	0.87	0.93	0.07
3	91.32%	0.85	0.91	0.0897
2	90.22%	0.82	0.9	0.092
1	89.72%	0.798	0.89	0.095

Table 10: Performance of prediction drawn by Stacked Bi-directional GRU model.

Number of Layer Used	Accuracy	Cohen-Kappa Score	F1-Score	MSE
4	84.37%	0.78	0.84	0.15
3	81.72%	0.76	0.82	0.178
2	80.6%	0.75	0.81	0.19
1	79.72%	0.73	0.8	0.298

As shown in Table 9 and Table 10, it is clear that the stacked bidirectional GRU model does not show promising efficiency as that of stacked bidirectional LSTM model. Hence, this model can be regarded as the best one for pursuing the CVD classification problem. Early prediction of heart disease may increase life span of the heart patient due to arised anxiety for numerous reasons. Considering past health record of a patient, the proposed Stacked Bidirectional LSTM Model can predict cardiac disease probabilities efficiently. This will assist the medical care units as well as accompany the doctors so that counter measures such as surgeries, medicines can be suggested. This proposed method reaches a promising and significant result that is dedicated towards heart disease prediction. Experimental results have shown prediction accuracy of 93.22%, F1-score of 0.93, kappa score of 0.87 with MSE of 0.07.

It is importantant to note that the following types of noises are considered in the ECG wave. It is shown in Figure 6.

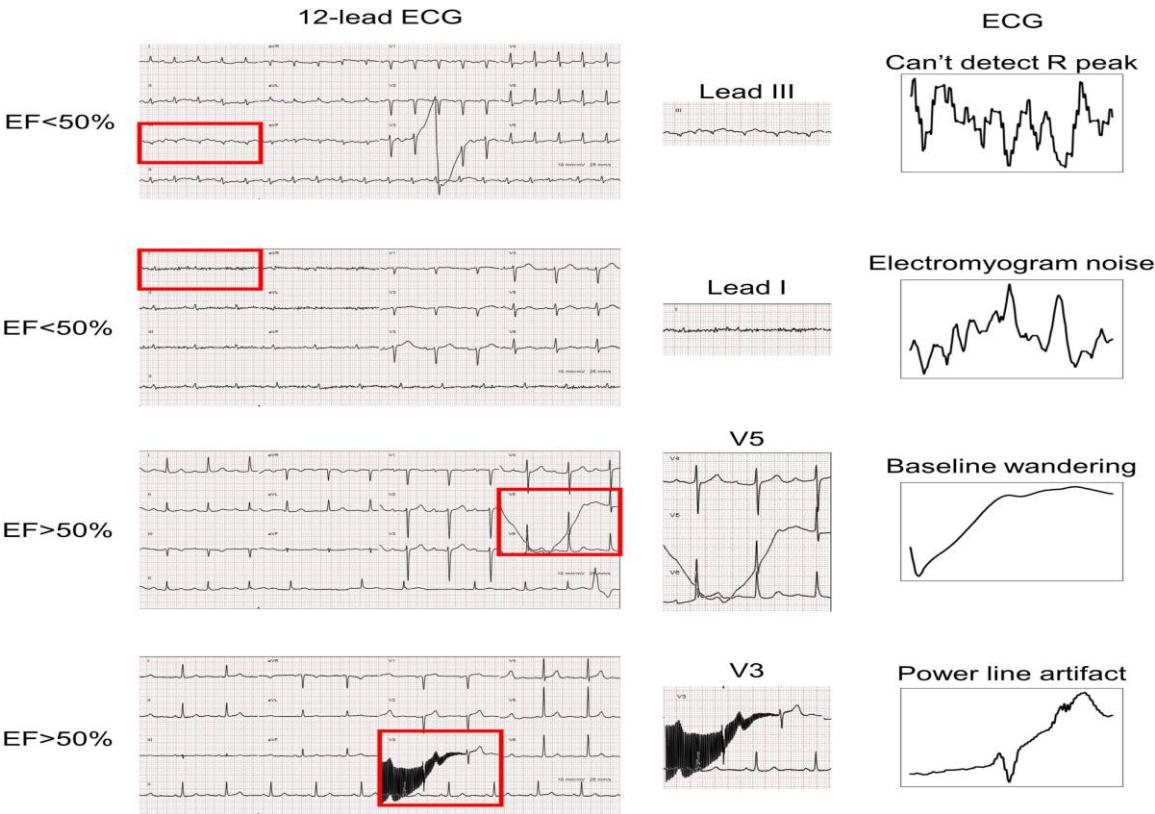


Figure 6 Various Types noises in ECG waves are considered

The following Table 11 will show types of cardiovascular diseases with symptoms, cause and prevention methods [39].

Table 11 Various Types of Cardiovascular Diseases

CVD Type	Symptoms	Cause	Prevention Methods
Heart Attack	Discomfort, Indigestion, Sweating, Vomiting, Irregular heartbeats.	Artery plaques attributable to calcium, fatty matter, proteins, and cells which are inflammatory.	Narcotics (aspirin, brilinta, etc.)surgical procedure Processes-Angioplasty
Coronary Heart Disease	Chest pain, Aching, Heaviness	Pulmonary embolism, Cardiomyopathy, Pericarditis,	Angioplasty, Bypass surgery.
Ischemic stroke	Headache, paralysis, or facial numbness, leg and arm, trouble with talking	Blocked artery hemorrhagic stroke.	Carotid endarterectomy, Angioplasty
Arrhythmia	Palpitations, fainting, dizziness, weakness, fatigue.	Electrolyte's incorrect balance in the blood, muscle changes in the heart	Medication, Change lifestyle, and surgery.
Heart valve Disease	Swelling of the feet, ankles, or abdomen, trouble with breathing and rapid gain in weight	Acquired valve disease, Congenital valve disease, Rheumatic fever	Medication, brush carefully to prevent teeth and gums infection
Enlarged Heart (Cardiomegaly)	Shortness of breath, weight gain, fatigue and leg swelling	Genetic and inherited conditions, infection of HIV, abnormal heart valve, high blood pressure.	Cardiac catheterization, high-blood regulation pressure, Avoiding the usage of harmful alcohol substances and caffeine
Heart Murmurs	High Blood Pressure and Anemia	Fever and hyperactive thyroid,.	Prevention of blood clots, surgery and diuretics through medicines
Cardiac Arrest	Racing Heartbeat, Dizziness	Abnormal Heart rhythms (Arrhythmia)	Consistently following-up with the doctors, surgery and medication

## Conclusions

Healthcare shows a significant key for perceiving the health related aspects of the humans around the globe. This chapter focuses on identifying CVDs in human heart from two perspectives. These two approaches cover syntactical pattern discover from ECG reports as well as construction of a predictive modeling using deep learning technique. The pattern discovery approach is an interesting domain yet challenging to perform because of its dependency on formal language generation. The predictive modeling is based on deep neural network. Multiple neural networks are utilized as the second approach. Use of neural network requires to be exemplified as it simulates human brain like tasks. Construction of an intelligent computerized tool is favored in this study as it facilitates the CVD classification task. Separating the CVD patients may assist the medical care unit to put more attention for their treatment. This task will definitely benefit the clinicians to assist in taking informed decisions.

## Conflicts of Interest

The authors declare that they have no conflicts of interest.

## References

1. Pagidipati, Neha Jadeja, and Thomas A. Gaziano. "Estimating deaths from cardiovascular disease: a review of global methodologies of mortality measurement." *Circulation* 127.6 (2013): 749-756.
2. Murray, Christopher JL, et al. "Using verbal autopsy to measure causes of death: the comparative performance of existing methods." *BMC medicine* 12.1 (2014): 1-19.
3. Buckberg, Gerald D., et al. "What is the heart? Anatomy, function, pathophysiology, and misconceptions." *Journal of cardiovascular development and disease* 5.2 (2018): 33.
4. Mall, Franklin P. "On the development of the human heart." *American Journal of Anatomy* 13.3 (1912): 249-298.
5. Rajkumar RP. COVID-19 and mental health: A review of the existing literature. *Asian journal of psychiatry*. 2020 Apr 10:102066.
6. Gonzalez, Rafael C., and Michael G. Thomason. "Syntactic pattern recognition: An introduction." (1978).
7. Stewart, Jack et al. "Primary prevention of cardiovascular disease: A review of contemporary guidance and literature." *JRSM cardiovascular disease* vol. 6 2048004016687211. 1 Jan. 2017, doi:10.1177/2048004016687211
8. Samir Kumar Bandyopadhyay, and Shawni Dutta. "Stacked Bi-directional LSTM Layer Based Model for Prediction of Possible Heart Disease during Lockdown Period of COVID-19: Bidirectional LSTM." *Journal of Advanced Research in Medical Science & Technology (ISSN: 2394-6539)* 7.2 (2020): 10-14.
9. Costa-jussà, Marta R., et al. "Introduction to the special issue on deep learning approaches for machine translation." *Computer Speech & Language* 46 (2017): 367-373.
10. Yu, Yong, et al. "A review of recurrent neural networks: LSTM cells and network architectures." *Neural computation* 31.7 (2019): 1235-1270.
11. Wang, Zhuo, et al. "Death burden of high systolic blood pressure in Sichuan Southwest China 1990–2030." *BMC public health* 20 (2020): 1-9.
12. Udupa, Jayaram K., and Ivaturi SN Murthy. "Syntactic approach to ECG rhythm analysis." *IEEE Transactions on Biomedical Engineering* 7 (1980): 370-375.
13. Albus, John Edward, et al. *Syntactic pattern recognition, applications*. Vol. 14. Springer Science & Business Media, 2012.
14. Goldschager, N., and Mervin J. Goldman. "Principles of clinical electrocardiography." (1989).
15. Cady, Lee D. *Computer Techniques In Cardiology*. New York: M. Dekker, 1979.
16. A. Chauhan, A. Jain, P. Sharma, and V. Deep, "Heart Disease Prediction using Evolutionary Rule Learning," *Int. Conf. Computational Intell. Commun. Technol. CICT 2018*, no. Cict, pp. 1–4, 2018, doi: 10.1109/CICT.2018.8480271.
17. A. H. Gonsalves, F. Thabtah, R. M. A. Mohammad, and G. Singh, "Prediction of coronary heart disease using machine learning: An experimental analysis," *ACM Int. Conf. Proceeding Ser.*, pp. 51–56, 2019, doi: 10.1145/3342999.3343015.



18. P. C. Austin, J. V. Tu, J. E. Ho, D. Levy, and D. S. Lee, "Using methods from the data-mining and machine-learning literature for disease classification and prediction: A case study examining classification of heart failure subtypes," *J. Clin. Epidemiol.*, vol. 66, no. 4, pp. 398–407, 2013, doi: 10.1016/j.jclinepi.2012.11.008.
19. M. Kirmani, "Cardiovascular Disease Prediction using Data Mining Techniques," *Orient. J. Comput. Sci. Technol.*, vol. 10, no. 2, pp. 520–528, 2017, doi: 10.13005/ojcs/10.02.38.
20. P. P. Sai and C. Reddy, "International Journal of Computer Science and Mobile Computing HEART DISEASE PREDICTION USING ANN ALGORITHM IN DATA MINING," *Int. J. Comput. Sci. Mob. Comput.*, vol. 6, no. 4, pp. 168–172, 2017.
21. B. Bahrami and M. Hosseini Shirvani, "Prediction and Diagnosis of Heart Disease by Data Mining Techniques," *J. Multidiscip. Eng. Sci. Technol.*, vol. 2, no. 2, pp. 3159–40, 2015.
22. S. Mohan, C. Thirumalai, and G. Srivastava, "Effective heart disease prediction using hybrid machine learning techniques," *IEEE Access*, vol. 7, pp. 81542–81554, 2019, doi: 10.1109/ACCESS.2019.2923707.
23. E. Izci, M. A. Ozdemir, M. Degirmenci, and A. Akan, "Cardiac arrhythmia detection from 2d ecg images by using deep learning technique," *TIPTEKNO 2019 - Tip Teknol. Kongresi*, pp. 1–4, 2019, doi: 10.1109/TIPTEKNO.2019.8895011.
24. C. Blake, E. Keogh, and C.J. Merz, "UCI repository of machine learning databases" [<http://www.ics.uci.edu/~mllearn/MLRepository.html>], Department of Information and Computer Science, University of California, Irvine, CA, 1998.
25. Vijayavanan, M., V. Rathikarani, and P. Dhanalakshmi. "Automatic classification of ECG signal for heart disease diagnosis using morphological features." *International Journal of Computer Science & Engineering Technology* 5.4 (2014): 449-455.
26. Karlik, Bekir, and A. Vehbi Olgac. "Performance analysis of various activation functions in generalized MLP architectures of neural networks." *International Journal of Artificial Intelligence and Expert Systems* 1.4 (2011): 111-122.
27. Srivastava, Nitish. "Improving neural networks with dropout." *University of Toronto* 182.566 (2013): 7.
28. Kline, Douglas M., and Victor L. Berardi. "Revisiting squared-error and cross-entropy functions for training neural network classifiers." *Neural Computing & Applications* 14.4 (2005): 310-318.
29. Brownlee, Jason. "What is the Difference Between a Batch and an Epoch in a Neural Network?." *Deep Learning; Machine Learning Mastery: Vermont, VIC, Australia* (2018).
30. Juba, Brendan, and Hai S. Le. "Precision-recall versus accuracy and the role of large data sets." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 33. No. 01. 2019.
31. Berry, Kenneth J., and Paul W. Mielke Jr. "A generalization of Cohen's kappa agreement measure to interval measurement and multiple raters." *Educational and Psychological Measurement* 48.4 (1988): 921-933.
32. Backus, John, 1978-Can programming be liberated from the von Neumann style? A functional style and its algebra of programs. *Communications of the ACM*, 21:613.
33. S. Hu, R. Gao, and L. Liu, "Summary of the 2018 report on cardiovascular diseases in China," *Chinese Journal of Circulation*, vol. 34, no. 3, pp. 209–220, 2019.

34. Y. Ping, C. Chen, L. Wu, Y. Wang, and M. Shu, "Automatic detection of atrial fibrillation based on CNN-LSTM and shortcut connection," *Healthcare*, vol. 8, no. 2, p. 139, 2020.
35. Q. Wu, Y. Sun, H. Yan, and X. Wu, "ECG signal classification with binarized convolutional neural network," *Computers in Biology and Medicine*, vol. 121, article 103800, 2020.
36. F. Ma, J. Zhang, W. Chen, W. Liang, and W. Yang, "An automatic system for atrial fibrillation by using a CNN-LSTM Model," *Discrete Dynamics in Nature and Society*, vol. 2020, Article ID 3198783, 9 pages, 2020.
37. A. K. Sangaiah, M. Arumugam, and G. B. Bian, "An intelligent learning approach for improving ECG signal classification and arrhythmia analysis," *Artificial Intelligence in Medicine*, vol. 103, article 101788, 2020.
38. J. Park, J. Kim, S. Jung, Y. Gil, J. I. Choi, and H. S. Son, "ECG signal multi-classification model based on squeeze-and-excitation residual neural Networks," *Applied Sciences*, vol. 10, no. 18, article 6495, 2020.
39. Felman, A. *Cardiovascular Disease: Types, Symptoms, Prevention, and Causes*. 2019.