



# MACHINE LEARNING ALGORITHMS IN PREDICTION OF HEART ATTACK

Jerushaa Jane K H S<sup>1</sup>, Dr.J.Kasthuri<sup>2</sup>, Aravindh S<sup>1</sup>, Dr.M.Buvana<sup>1\*</sup>, Dharshan Kumaar S<sup>1</sup> and Kailash U V<sup>1</sup>

<sup>1</sup>Department of Computer Science and Engineering

<sup>1\*</sup>Associate Professor, <sup>1</sup>Department of Computer Science and Engineering  
PSNA College of Engineering and Technology, Dindigul- 624005, Tamil Nadu, India

<sup>2</sup>Assistant Professor, Department of Zoology,

The Standard Fireworks Rajaratnam College, Sivakasi- 626124, Tamil Nadu, India

**Abstract:** In the global scenario, prevention of heart attack has become imperative in decreasing the death rates. The chance of occurrence of heart attack can be predicted with proper dataset. With the rapid advancement in Information Technology data pertaining to health care issues are being updated on a daily basis. The collected data are then authenticated by analyzing them with various machine learning algorithms. As the cardio professionals have their own limitations to predict the chance of getting heart attack to high accuracy, the present investigation attempts to process and interpret the dataset of prominent heart attack patients. Further, by training the chosen models, predictions are being made with about 6 machine learning algorithms. Results of the various machine learning algorithms reveal that Kaggle recorded dataset entitled "Heart attack prediction" is seemed to be working more efficiently with Logistic Regression and showing the highest accuracy (91.8%).

**Index Terms** - Heart attack Predictions, Kaggle Dataset, Machine Learning Algorithms, Logistic Regression.

## I INTRODUCTION

'Heart Attack' (HA) or myocardial infarction is a life threatening experience. It normally occurs when there is a restriction in flow of blood to some part of the heart muscles. It often results from a blockage in a nearby artery. This blockage may be caused by the accumulation of fat, cholesterol and some other substances which could form plaques / clots in the coronary arteries. The condition of developing plaques in coronary artery is called atherosclerosis. However, development of plaques takes several years. During the course of time, plaques may either become hardened or ruptured. Hardened plaques eventually narrow down the coronary arteries and thereby reduce the flow of oxygenated blood to heart. Besides, on rupturing of plaques, blood clots are being formed on its surface. Evidently, a larger sized blood clot terribly blocks the blood flow. Over time, the ruptured plaque also hardens and narrow downs the coronary arteries. Thus the partially blocked / suppressed blood flow caused by plaque / rupture in turn impairs or destroys a portion of the heart muscle. Without proper diagnosis and therapy in time HA becomes fatal. Remarkably, HA turns into a common cause of death worldwide.

Alarmingly, in United States a death occurs for every 40 seconds due to cardiovascular disease (CVD). Thus in a year, about 6,55,000 Americans die. According to the recent studies, in India on an average, 50% population at the age group of 50 plus years and 25% of 40 plus years are experiencing HA. Especially, people who live in cities are three times more vulnerable to HA than those of them living in villages. It is noteworthy that, many people do survive with HA due to persistent treatment for over a period of years. In the present scenario, spotting the early signs / symptoms and taking prompt treatment can save life. Thus realizing the need for accurate prediction of HA in a short time, the present investigation attempts to incorporate some of the machine learning (ML) techniques to analyze the historical data. The algorithms are then adopted by computers for computation and solving the problems. Machine learning algorithms makes a difference by allowing the computers to get trained on input datasets and to make use of statistical methods so as to analyze the output (results). These supervised learning approaches also facilitate both organizations and people to tackle real world problems such as identifying the malware infected links and discovering the hidden malwares in a file using steganographic techniques. Surprisingly, the file with hidden malwares seems to appear as normal. Hence with the view to improve human lives the present investigation has been designed to work out with the following models of ML in HA prediction.

- Logistic Regression
- Naïve Bayes
- Decision Tree
- Random Forest
- Support Vector Machine
- K-Nearest Neighbor

Further, by visualizing and analyzing the results with the chosen dataset, the one which could show highest accuracy shall be considered as the best suited ML algorithm for HA prediction.

## II. LITERATURE SURVEY

In the global scenario HA takes the life of about 12 million people every year. For its high prevalence, it needs to be diagnosed timely and effectively (Ramalingam *et al.*, 2018). As diagnosis / prediction of heart disease (HD) is a great skill requiring complex task, medical organizations, all around the world collect data on various health related issues pertaining to cardiovascular diseases (CVD). However, these datasets are very massive and remains noisy. Besides, they are too vast for human minds to comprehend. Nevertheless, they can be easily explored by using various machine learning techniques to gain useful insights (Vapnik, 1995). In recent times, ML algorithms have started paving an easy way for a win-win situation in accurate prediction of health issues with special reference to HA / HD. These HD prediction systems enable the health care practitioners to predict the status of HA based on the clinical data of patients. In this regard, the scientists do evaluate datasets with ML algorithms and provide easy output. Thus the process of binding the records of medical data together digs novel techniques and enables us to evaluate results in a new dimension.

Accordingly, Mrudula *et al.* (2010) have suggested Support Vector Machine (SVM) and Artificial Neural Network (ANN) methods as the decision support systems for classification of HD. Especially for ANN they have chosen a Multilayer Perceptron Neural Network (MLPNN) to develop a decision support system. This MLPNN has been trained by back – propagation algorithm and proved to be a computationally efficient method for diagnosing HD. Whereas, Kavitha *et al.* (2010) have designed an evolutionary neural network for the detection of HD. Asha and Sophia Reena (2010) on the other hand have worked out with the supervised machine learning based classification model for diagnosis of HD. Chen *et al.* (2011) have presented a computational model for HD prediction system using artificial neural network algorithm with the data comprising 13 important clinical features *viz.*, age, gender, chest pain type etc., collected from machine learning repository of UCI. Further, Cinetha and Uma Maheswari (2014), have adopted fuzzy logic as a decision support system to rule out coronary heart disease. Their system successfully predicts the possibility of incidence of HD in a person for the next ten years.

Furthermore, Sairabi and Deval (2015) have predicted HD with modified k-means and Naïve Bayes. Suganya and Tamije Selvy (2016) for their turn implemented Fuzzy Cart Algorithm for HD prediction. In this case fuzziness is seemed to be introduced in the measured data to eliminate uncertainty. While Ashwini Shetty and Chandra (2016) have adopted neural network and genetic algorithm to predict HD, Manpreet *et al.* (2016) have proposed the Structural Equation Modeling (SEM) and Fuzzy Cognitive Map (FCM) based HD prediction system. For which they have chosen twenty significant attributes from Canadian Community Health Survey dataset. Where SEM generates weight matrix and FCM predicts the possibility of cardiovascular disease (CVD). Subsequently, Prajakta *et al.* (2015) have developed a smart HA prediction system by harnessing big data and data mining modeling techniques.

In a yet another study, Sharmila and Indra Gandhi, (2017) have proposed a conceptual method for the prediction of HD using data mining techniques. Similarly, Ramalingam *et al.* (2018) have exploited the ML algorithms on complex medical datasets. Pushpalatha (2019) on the other hand has proposed data mining techniques for the analysis of HD. Data mining is the exploration of large datasets. It extracts hidden and previously unknown patterns, relationships and knowledge that are difficult to detect with traditional statistical methods (Liao *et al.*, 2000). All the more, it is worth mentioning that, the ML algorithms are not only employed for accurate prediction of HD but also to suggest the medical prescription. Besides, the supervised machine learning method can also be employed to classify the patients (Dharshana Deepthi *et al.*, 2020). In this regard, Shaik *et al.* (2010) have adopted Radial Basis Function. As a consequence the present investigation attempts to exploit some of the supervised ML algorithms for effective analysis of different categories of chosen datasets (age, sex, cp, trestbps, cholesterol, fbs, restecg, thalach, exang, oldpeak, slope, ca, thal & target) showing major differences and to extract the results with highest accuracy levels.

## III. DATASET

In the present investigation, the dataset chosen (Fig. 1) for “Heart Attack Prediction” has been accessed from Kaggle (Online dataset). It is worth mentioning that, this dataset is also being adopted by World Health Organization (WHO).

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1

Figure 1:Dataset chosen from Kaggle

Thus adoption of this dataset has become imperative in a pre-processing study pertaining to HA prediction. Description of the features of a dataset (Table 1), of the present investigation is as follows.

**Table 1: Description of the Dataset**

S.No.	Feature	Description	Range
1	Age	Age	29 to 77
2	Sex	Gender	0 and 1
3	Chol	Cholesterol	126 to 564
4	Fbs	Fasting Blood Sugar	0 to 1
5	Thalach	Maximum Heart-Rate	71 to 202
6	Restecg	Resting ECG	0 to 2
7	Oldpeak	ST depression by exercise	0 to 6.2
8	Trestbps	Resting blood pressure	94 to 200
9	Slope	Slope of peak exercise ST	0 to 2
10	Cp	Chest pain type	0 to 2
11	Ca	Major Vessels	0 to 4
12	Thal	Heart defect	0 to 3
13	Exang	Angina due to exercise	0 to 1

## IV. ARCHITECTURAL FRAMEWORK

### 4.1. Feature Extraction

Feature extraction is a practice of identifying and combining some of the major attributes / properties of a dataset to be employed for working out the model. In other words, it is a process of finding out the subset of original features so as to achieve certain goals. Features / properties of the dataset chosen for prediction of HA in the present investigation are distinguishable from one another. They not only improve accuracy but also reduce the processing time.

### 4.2. Instance Labeling

Nowadays data labeling is playing crucial role in producing an efficient data model. Particularly, ML is making use of enormous amount of data for working out the model. In the present investigation, data used in ML are annotated / labeled and organized. Hence, the model would make use of the organized data to predict HA more accurately and effectively.

### 4.3. Training the Model

Machine Learning Model is being guided to make use of the training dataset which contributes 80% of data in the dataset. This approach of training the model makes prediction more accurate than ever before. Further, as both the quantity and quality of dataset equally play significant role in HA prediction with high accuracy, the training data are required to be more precisely annotated and cleaned.

### 4.4. Data Pre-Processing and Cleaning

Dataset chosen for the present investigation contains variables of different data types. It also includes some categorical variables. Label encoding practice makes use of absolute values of dataset and labels the distinct absolute values. The problems encountered due to the missing values of variables, can be rectified by assigning "NA" to them. In the present study, data pertaining to death of HA patients, recovery and missing values are compiled for testing. Besides, the well defined records are also being compiled for training (Fig. 2). Further, the complex data are formatted to achieve better results. In this context, the transformation process has been adopted to change the data format from one form to other and making them most comprehensible. This is being achieved through the techniques *viz.*, normalization, smoothing, generalization and aggregation of data.

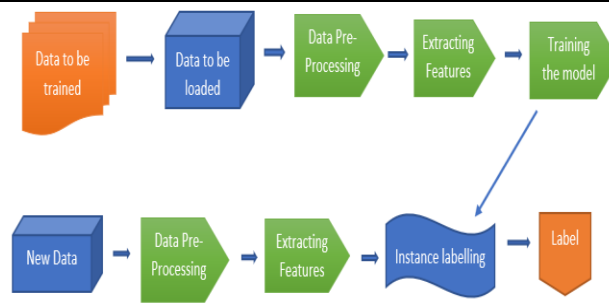
```

from sklearn.model_selection import train_test_split
y=df['target']
x=df.drop('target',axis=1)
x_train,x_test,y_train,y_test=train_test_split(x,y,train_size=0.8,test_size=0.2)
x_train.head()

```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
90	48	1	2	124	255	1	1	175	0	0.0	2	2	2
137	62	1	1	128	208	1	0	140	0	0.0	2	0	2
115	37	0	2	120	215	0	1	170	0	0.0	2	0	2
127	67	0	2	152	277	0	1	172	0	0.0	2	1	2
73	51	1	0	140	261	0	0	186	1	0.0	2	0	2

**Figure 2: Dataset after training and testing**



**Figure 3: Architectural Diagram of the proposed model**

Further, for easy understanding of the principles behind the working of algorithms, architectural diagram has also been proposed (Fig. 3).

## V. MACHINE LEARNING ALGORITHMS

### 5.1. Logistic Regression

Logistic regression (LR), a classification algorithm (Sperandei, 2014) for categorical variables is being used for the prediction of probability of occurrence of chosen target variables. The target variable is predicted in such a way to show only two possible outcomes. The outcomes in turn are encoded as '1' to indicate success and '0' for failure. This is one of the simplest ML algorithms. It can also be used for classification of various problems *viz.*, blood pressure prediction, spam detection, diabetes prediction, relationship identification between micro RNA and gene, credit scoring, etc. This Logistic Regression makes use of 'Sigmoid function' and configures its resultant value to be either 1 or 0. For better prediction, sigmoid function is being used for effective mapping of probabilities between 0 and 1. Further, better prediction can also be achieved by avoiding multi co-linearity and meaningless variables in the large chosen dataset.

$$f(x) = 1/(1 + e^{-x})$$

### 5.2. Gaussian Naïve Bayes

Extremely huge datasets can be easily interpreted and analyzed using Gaussian Naïve Bayes (GNB) algorithm of supervised machine learning technique (Zang, 2005). Recommendation systems using machine learning and data mining can be effectively built with Naïve Bayes and treated for collaborative filtering for prediction / filtration of unseen information. This algorithm makes use of Bayes theorem which explains the event's probability based on the prerequisite of event. Naïve Bayes assumes that all features of an event are independent of other and whose presence or absence does not influence the other features.

$$P(H | E) = (P(E | H) * P(H)) / P(E)$$

Where P(H) is the probability hypothesis of being true. It is also known as Prior or Previous Probability. Whereas P(E) is the Evidence Probability regardless of hypothesis. P(H | E) is the Hypothesis Probability indicating the occurrence of event. P(E | H) is the evidence probability indicating that the hypothesis is true.

### 5.3. Decision Tree

The most popular tool (Mitchell, 1997) in ML, commonly interpreted as a decision support tool is making use of a tree like graph or model of decisions and their possible consequences including chance, event outcomes and utility. This structured algorithm is capable of making efficient decision. A decision tree can be easily transformed to a set of classification rules by mapping with its path from root node to leaf node (Top down Approach). In this case, competent decision shall be taken at every node based on the train and test data chosen for efficient traversal of the tree. This tool simplifies a complex logic by classifying the decisions. Clients choose decision tree algorithm over any algorithms for its simplicity, explanation and usefulness. Generally these decision trees are being used in research operations, especially in decision analysis to help and identify a strategy that will most likely to lead to reach the goal.

### 5.4. Random Forest

Random Forest (RF) an ensemble classifier (Breiman, 2001), remains popular among the supervised machine learning techniques. It is being adopted in classification and regression problems for its wide applications in remote sensing, object detection and Kinect gaming console. Even then there is a large data missing, it is being recommended for its high accuracy in prediction. It makes use of multiple decision trees with less training time to decrease the risk of over fitting. Predictive analysis is being made out from an average of the trees. It works by combining many classifiers to solve complex problems and to increase the performance. With this algorithm, accuracy can also be predicted even when a large amount of missing data stands as a barrier. Further, it minimizes the problems of high variance and high bias by averaging to find a natural balance between the two extremes.

### 5.5. Support Vector Machine

Support vector machine (SVM) is being categorized as supervised machine learning method (Ben-Hur *et al.*2001). It has been used for classification and regression enigmas. SVM creates decision boundary which in turn separate n-dimensional space into classes so as to position the new upcoming data into a place or point related to its category. The decision boundary is also known as hyperplane. It can be created by finding out the correct extreme vector points. Data points or vectors, closest to hyperplane and the one could affect the position of hyperplane are termed as Support Vector. There can be many decision boundaries for segregation of classes in n-dimensional space. As a rule an efficient data boundary should classify the data points. The dataset features in turn should determine the hyperplane dimension. They do have maximum margins in hyperplane and are being used for the indication of data points with maximum distance. To a core, SVM attempts to maximize the margin (distance between hyperplane and the two closest data points) to decrease the chance of misclassification

### 5.6. K-Nearest Neighbor Algorithm

K-Nearest Neighbor algorithm of supervised machine learning is a nonparametric technique for pattern classification (Fix and Hodges, 1951). It assumes similarities between the new case data and available data. It can be adopted not only in classification but also in regression analysis. Actually, K-NN algorithm stores the available data and classifies the new one based on their similarity. It groups the new case data in accordance with the most similar available category. Being non-parametric algorithm it does not make any assumptions with underlying data. It is also called lazy learner algorithm as it does not learn from the training set, instead it stores dataset during classification and perform action on it. K-NN algorithm uses similarity of the feature for predicting values of new data points and assign them with values on analyzing how closely the data points match with training set points. The following steps have to be taken into account in K-NN algorithm:

- Loading of train data and test data into a code.
- Selection of K value nearest to the data point.
- Adoption of Euclidean, Manhattan or Hamming distance methodologies to calculate the span of test data with its corresponding train data.
- Sorting the calculated distances in descending order and selecting K rows from the top.
- Assigning a test data point with a class by analyzing the frequent classes in a row.

## VI. WORK FLOW

The present investigation on HA prediction has been meticulously designed to execute with the following work flow (Fig. 4) model.

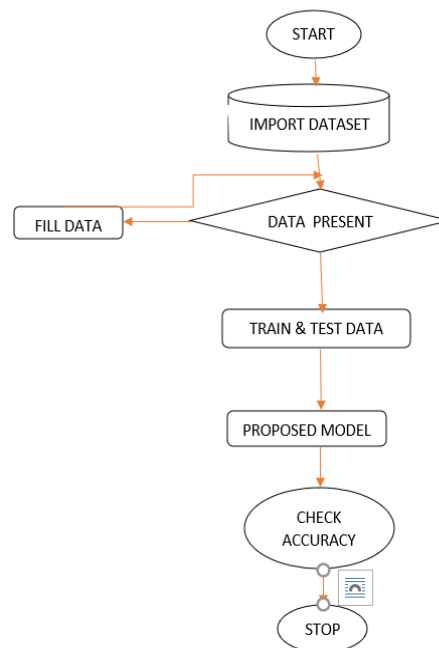


Figure 4: Work Flow

## VII. EVALUATIONARY STUDY

Evaluation study makes use of True Positive (TP), True Negative (TN), False Positive (FP) and False Negative (FN) values. The corresponding data points of TP, TN, FP and FN obtained against the chosen algorithms of the present study are being presented (Table 2) after making a serious assessment.

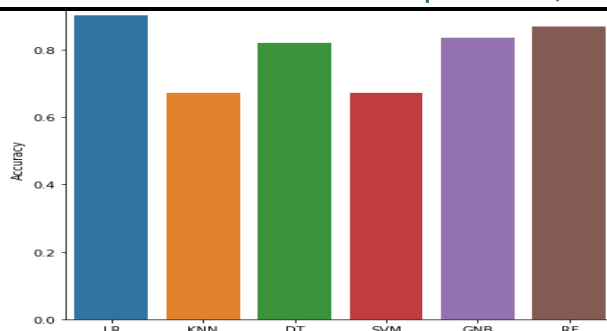
Table 2: Data points of select algorithms

ML Algorithms	Data Points			
	TP	TN	FP	FN
Logistic Regression	32	19	7	3
Gaussian NB	33	19	7	2
Decision tree	28	18	8	7
Random Forest	32	20	6	3
SVM	27	18	8	8
KNN	21	22	4	14

### 7.1. Accuracy

The dataset contains True Positive (TP) and True Negative (TN) points. Accuracy of dataset has been calculated as the ratio of sum of data points of TP and TN to the sum of TP, TN, FP and FN data points as in the Formula (1). Accuracy estimates the performance of the algorithmic model. It lies between  $0.0 <$  and  $< 1.0$ . Accuracy for the present investigation is being best described in Figure 5.

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} \quad (1)$$

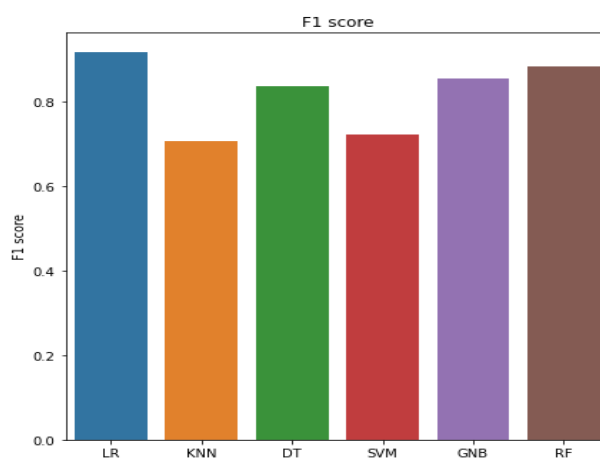


**Figure 5: Accuracy**

### 7.2. F1 Score

F1 score (Formula 2) is highly essential when processing an uneven class distribution providing harmony between recall and precision as illustrated (Fig.6). It is calculated as twice the ratio of the product of precision and recall to the sum of precision and recall.

$$\text{F1 Score} = 2 * \frac{\text{Precision} * \text{recall}}{\text{Precision} + \text{recall}} \quad (2)$$

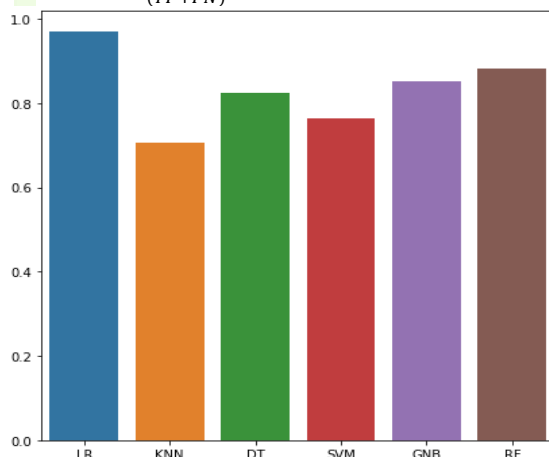


**Figure 6: F1 score**

### 7.3. Recall

Recall defines the number of true positives recorded. It is the number of correctly classified patients in an imbalanced class dataset, out of all the patients who have been correctly predicted for HA (Fig. 7). It can be worked out with a formula (3).

$$\text{Recall} = \frac{(TP)}{(TP+FN)} \quad (3)$$



**Figure 7 : Recall**

### 7.4. Precision

Precision evaluates (formula 4) whether the labeled true positives are really true positives. It fabricates a value that determines correctly classified HA patients (Fig. 8). It is worth mentioning that high precision relates to low false positive rates.

$$\text{Precision} = \frac{TP}{TP+FP} \quad (4)$$

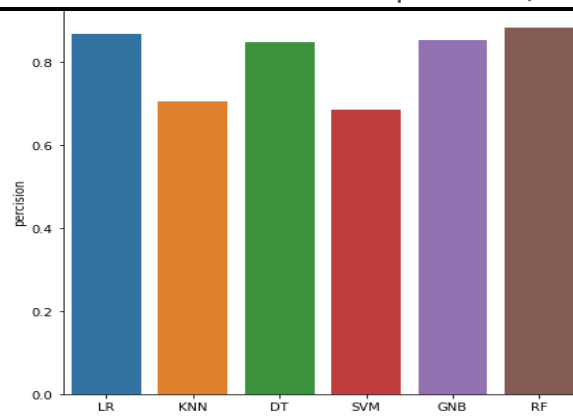


Figure 8: Precision

VIII. RESULTS AND DISCUSSION

Findings of the present investigation on predictive analysis of HA (Table 3) emphasize that the Logistic Regression confers the highest chance of acquiring HA in patients based on different parameters with 0.918033 Accuracy, 0.875000 Precision, 0.965517 Recall and 0.918033 F1 score. These values are comparatively greater than K-Nearest Neighbor (Accuracy: 0.672, Precision: 0.63, Recall: 0.724, F1 score: 0.67), Decision tree (Accuracy:0.81, Precision: 0.80, Recall: 0.82, F1 score: 0.81), Support Vector Machine (Accuracy: 0.65, Precision: 0.61, Recall: 0.724, F1 score: 0.66), Gaussian NB (Accuracy: 0.86 , Precision: 0.81, Recall: 0.93, F1 score: 0.87) and Random Forest (Accuracy: 0.86, Precision: 0.81, Recall: 0.93, F1 score: 0.87) algorithms as interpreted from Figure 9 and 10.

```
comparison = pd.DataFrame({
    "Logistic regression":{'Accuracy':log_accuracy, 'percision':logistic_percision,'recall':logistic_recall,'F1 score':logistic_f1},
    "K-nearest neighbours":{'Accuracy':knn_accuracy, 'percision':knn_percision,'recall':knn_recall,'F1 score':knn_f1},
    "Decision trees":{'Accuracy':tree_accuracy, 'percision':tree_percision,'recall':tree_recall,'F1 score':tree_f1},
    "Support vector machine":{'Accuracy':svm_accuracy, 'percision':svm_percision,'recall':svm_recall,'F1 score':svm_f1},
    "Gaussian NB":{'Accuracy':clf_accuracy, 'percision':clf_percision,'recall':clf_recall,'F1 score':clf_f1},
    "Random Forest":{'Accuracy':rfc_accuracy, 'percision':rfc_percision,'recall':rfc_recall,'F1 score':rfc_f1}
})
comparison.head()
```

Figure 9 : Algorithms vs features

Table 3: Performance of various algorithms over the dataset

Features	Machine Learning Algorithms					
	Logistic Regression	Gaussian NB	Decision Tree	Random Forest	Support Vector Machine	K-Nearest Neighbours
Accuracy	0.918033	0.868852	0.819672	0.868852	0.655738	0.672131
Precision	0.875000	0.818182	0.800000	0.818182	0.617647	0.636364
Recall	0.965517	0.931034	0.827586	0.931034	0.724138	0.724138
F1Score	0.918033	0.870968	0.813559	0.870968	0.666667	0.677419

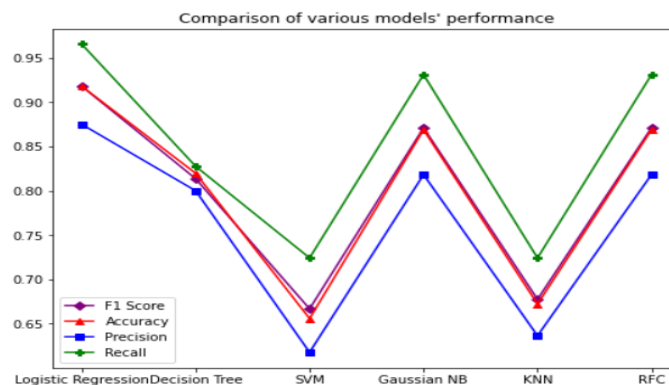


Figure 10: Comparison of model's performance

Thus HA prediction analysis of the chosen dataset treated with algorithmic techniques reveals high accuracy, precision, recall and F1 score for Logistic Regression Algorithm. Whereas, Dwivedi (2018) while working out with certain algorithms of machine Language viz., navies bayes, KNN and logistic regression for accurate prediction of HD against the given datasets, results of their findings showed highest accuracy with NB when compared to the other algorithms. Xu *et al.* (2017), on the other hand, while treating Cleveland dataset with various methods, have reported the highest accuracy (91.6%) against random forest. Further, the conceptual model of HD prediction proposed by Sharmila and Indra Gandhi, (2017), reveals better and efficient accuracy (85%) with SVM.

## IX. CONCLUSION

Findings of the present investigation on HA prediction analysis with various machine learning algorithms reveal that the dataset from Kaggle entitled “Heart attack prediction” is seemed to be working more efficiently with Logistic Regression and showing 91.8% accuracy. This value remains relatively high when compared to K-NN, Decision tree, SVM, Gaussian NB and Random Forest Algorithms. Hence, the authors arrive to a conclusion that the Logistic Regression as the one of the best suited algorithm for HA prediction. This algorithm can be implemented by the researchers / developers as the solution for HA prediction.

## X. REFERENCES

- [1] Asha, R. and Sophia Reena, G. 2010. Diagnosis of Heart Disease using Data Mining Algorithms. *Global Journal of Computer Science and Technology*, 10 (10): 38-43.
- [2] Ashwini Shetty, A. and Chandra, N. 2016. Different Data Mining Approaches for Predicting Heart Disease. *International Journal of Innovative Research in Science, Engineering and Technology*, 5(9) : 277-281.
- [3] Ben-Hur,A., Horn,D., Stegelmann, H.T. and Vapnik, V. 2001. Support Vector Clustering. *Journal of Machine Learning Research*, 2:125-137.
- [4] Breiman, L. 2001. Random Forests. *Machine Learning*, 45 : 5-32.
- [5] Chen, A.H., Huang, S.Y., Hong, P.S., Cheng, C.H. and Lin, E.J. 2011. HDPS: Heart Disease Prediction System. *Computing in Cardiology*, 38 : 557- 560.
- [6] Cinetha, K. and Uma Maheswari, P. 2014. Decision Support System for Precluding Coronary Heart Disease using Fuzzy Logic. *International Journal of Computer Science Trends and Technology*, 2 (2) : 102-107.
- [7] Dharshana Deepthi, L., Shanthi, D. and Buvana, M. 2020. An Intelligent Alzheimer’s Disease Prediction Using Convolutional Neural Network (CNN). *International Journal of Advanced Research in Engineering and Technology* 11(4):12-22.
- [8] Dwivedi, A. K. 2018. Performance evaluation of different machine learning techniques for prediction of heart disease. *Neural Computing and Applications*, 29(10)DOI:10.1007/s00521-016-2604-1.
- [9] Fix, E. and Hodges,J.L. 1951. Discriminatory analysis, nonparametric discrimination: Consistency properties. Technical Report 4, USAF School of Aviation Medicine, Randolph Field, Texas.
- [10]Kaggle [Online]Dataset: <https://www.kaggle.com/ronitf/heart-disease-uci>
- [11]Kavitha,K.S., Ramakrishnan, K.V. and Manoj K.S. 2010. Modeling and Design of Evolutionary Neural Network for Heart Disease Detection. *International Journal of Computer Science*, 7(5) :272-283.
- [12]Liao, W.B., Liu, C.F., Chiang,C.W., Kung,C.T. and Lee, C.W. 2000.Cardiovascular manifestations of pheochromocytoma. *American Journal of Emerging Medicine*, 18(5): 622-625.
- [13]Manpreet ,S., Levi M.M., Patrick J. and Vijay,K. M. 2016. Building a Cardiovascular Disease Predictive Model using Structural Equation Model and Fuzzy Cognitive Map. *IEEE International Conference on Fuzzy Systems (FUZZ)*, pp. 1377-1382.
- [14]Mitchell,T. 1997. *Machine Learning*. The McGraw-Hill Companies, Inc., pp. 52-78.
- [15]Mrudula, G., Kapil, W. and Snehlata, D. 2010. Decision Support System for Heart Disease Based on Support Vector Machine and Artificial Neural Network. *International Conference on Computer and Communication Technology*, DOI:10.1109/ICCCT.2010.5640377, 17-19.
- [16]Prajakta,G., Vrushali, G., Kajal, K. and Prajakta, D. 2015. Intelligent Heart Attack Prediction System Using Big Data. *International Journal of Recent Research in Mathematics Computer Science and Information Technology*, 2(2) : 73-77.
- [17]Pushpalatha, K. 2019. Analysis of Heart Disease Prediction System using Data Mining Techniques. *International Journal for Research in Applied Science and Engineering Technology*, 7(4): 8-14.
- [18]Ramalingam, V. V., Dandapath, A. and Karthik Raja, M. 2018. Heart disease prediction using machine learning techniques: A survey. *International Journal of Engineering and Technology*, 7: 684 -687.
- [19]Sairabi, H. M. and Devale, P. R. 2015.Prediction of Heart Disease using Modified k-means and by using Naive Bayes. *International Journal of Innovative Research in Computer and Communication Engineering*, 3(10) : 10265-10273.
- [20]Shaikh A.H., Mane, A.V., Manza, R. R. and Ramteke, R. J. 2010. Prediction of Heart Disease Medical Prescription using Radial Basis Function. *IEEE International Conference on Computational Intelligence and Computing Research (ICCIC)*, DOI: 10.1109/ICCIC.2010.5705900 ,28-29.
- [21]Xu,S., Zhang,Z., Wang,D., Hu,J., Duan,X. and Zhu,T. 2017. Cardiovascular Risk Prediction Method Based on CFS Subset Evaluation and Random Forest Classification Framework. *IEEE 2nd International Conference on Big Data Analysis (ICBDA)*,pp228-232. DOI: 10.1109/ICBDA2017.8078813.
- [22]Sharmila, S. and Indra Gandhi, M.P. 2017. Analysis of Cardiovascular Disease Prediction using Data Mining Techniques. *International Journal of Modern Computer Science*, 8 (5): 93-95.
- [23]Suganya, S. and Tamije Selvy, P. 2016. A Proficient Heart Disease Prediction Method using Fuzzy-Cart Algorithm. *International Journal of Scientific Engineering and Applied Science*, 2 (1): 1-6.
- [24]Sperandei, S. 2014. Understanding logistic regression analysis. *Biochem. Med.*, 24 (1): 12- 18.
- [25]Vapnik,V. 1995. *The Nature of Statistical Learning Theory*. Springer - Verlag, New York, pp.188.
- [26]Zhang, H. 2005. Exploring conditions for the optimality of .*International Journal of Pattern Recognition and Artificial Intelligence*, 19(2) : 183-198.