# FEATURE-BASED IMPLEMENTATION OF MACHINE LEARNING ALGORITHMS FOR CARDIOVASCULAR DISEASE PREDICTION

## H. Singh[1], R. Tripathy[2], P. Kumar Sarangi[3], U. Giri[4], S. Kumar Mohapatra[5], N. Rameshbhai Amin[6]

[1, 2, 3, 4] Chitkara University School of Engineering & Technology
Chitkara University, Himachal Pradesh, India.

[5]Chitkara University Institute of Engineering & Technology
Chitkara University, Punjab, India.

[6]Electronics & Communication Engineering Department, Government Engineering College, Bhavnagar, India.

Email: [1]hakam.singh@chitkarauniversity.edu.in
[2]ramamani.tripathy@chitkarauniversity.edu.in
[3]pradeepta.sarangi@chitkarauniversity.edu.in, [4]uttameast@gmail.com
[5]srikanta.mohapatra@chitkara.edu.in, [6]niravamin@gecbhavnagar.ac.in

Corresponding Author: **P. Kumar Sarangi**

## Abstract

*In eukaryotic organisms, each and every organ takes a major role in ensuring the seamless functioning of the entire system. If we consider about heart then it is treated as a vital part of every human being. Heart-associated ailments are very frequent at present so it is essential to predict such illnesses. This prognosis and prediction of coronary heart-associated illnesses require a lot of accuracy so it must be finished in an environment-friendly manner due to the fact a small mistake can motivate the death of the person. To deal with this hassle there ought to be a gadget which can predict and create consciousness about diseases. It is challenging to decide the ailment manually primarily based on signs and hazard factors. But this ought to be solved with the use of Machine mastering techniques. Artificial brain (AI) in the shape of desktop studying (ML) allows software program purposes to predict results greater precisely whilst functioning unbiased of human input. This study employs various machine learning algorithms, including K-Nearest Neighbors, Support Vector Machine, Logistic Regression, Random Forest, Decision Tree, and Naïve Bayes, to assess their accuracy in predicting cardiovascular disease and related conditions This paper makes use of the UCI repository dataset for coaching and testing including some basic parameters such as age and sex. After applying all algorithms to our data set, the experimental results concluded that the Logistic Regression model has predicted well with highest accuracy of 92% in comparison with other algorithms.*

*H. Singh et al.*

## I.  Introduction

The heart is the essential organ of your circulatory gadget. It pumps blood to every phase of our body. If there is any trouble in the functionality, then the intelligence and a range of different organs will quit working, and within a few minutes, the character will die [I].  In several medical applications nowadays, machine learning models are a widely used approach. [II] Cardiovascular ailments or heart disease are steadily growing day by day in our modern society. The World Health Organization (WHO) reports that the major factor for the cause of death globally is cardiovascular disease, which accounts for 17.9 million deaths yearly, or 32% of all deaths, including heart attacks and strokes. Almost 24.8% of deaths in India, according to the Global Burden of Disease, are caused by CVDs. [III] Identification of a coronary heart disorder is tough due to the fact there are countless dangerous elements like excessive blood sugar stage or diabetes, hypertension, excessive quantity of low-density lipoprotein (LDL) cholesterol, Uncontrolled stress, depression, and anger, and many extra factors. Many unique assessments like Electrocardiogram (ECG or EKG), exercise exams or stress tests, and Cardiac catheterization are carried out for the prognosis of cardiovascular diseases (CVDs).

Heart ailment prediction the use of computing device mastering techniques has been an ongoing effort. By taking into consideration elements like chest discomfort, ideal cholesterol levels, a person's age, and different factors, computing device studying may additionally be used to decide whether or not an individual has a cardiovascular disease. Machine mastering has the doable to considerably enhance the accuracy of coronary heart ailment prediction through inspecting giant quantities of scientific facts and figuring out patterns that might also point out a cardiac condition's existence.

This paper makes use of four supervised algorithms: Naïve Bayes, Support Vector Machine, Logistic Regression, and K-Nearest Neighbor. The elements and clinical records of the affected person are discovered in a dataset chosen from the UCI repository [IV] and after that different ML algorithms were implemented on this data set.

## II.     Related work

This research is inspired by several previous efforts in the literature on machine learning-based heart disease detection. This section highlights some notable works done by researchers before this work. To predict cardiac illnesses, authors like Patel et al. [V] implemented the Logistic Model Tree Algorithm and Random Forest Algorithms in their study. Using the data, the authors found that heart disease can be predicted accurately using the J48 tree technique. Soni et al. [VI], in their work, have used a Weighted Associative Classifier for heart attack prediction. The authors have developed GUI to input the patient's information, and using the rules saved in the rule base, the authors conclude that it was possible to forecast whether the patient had heart disease or not.

*H. Singh et al.*

Shao et al. [VII] in their work have implemented Naïve Bayes and Random Forest classifiers for forecasting using diversified datasets. According to the results produced by the authors, for some databases Random Forest performs better, and for others, Naïve Bayes performs better. Authors like Chauraisa et al. [VIII] have carried out research on respective data mining techniques to detect internal organ disease. The accuracy achieved by the authors in the bagging approach is 85.03% after using the three classifier algorithms to diagnose heart disease patients, and it took 0.05 seconds to create the model in total. Mrs. G. Subba Lakshmi al. [IX] the Heart Disease Prediction System's Decision Support was crafted utilizing the Naive Bayesian Classification Method The system searches a historical database of cardiac illnesses to find hidden data.

In another work, authors like Nikhar et al. [X] have implemented a Decision Tree classifier and Naïve Bayes for the prediction of heart diseases. From the experiment, it is revealed that Decision trees are more accurate than naive Bayes classifiers. To anticipate and diagnose cardiac disease, authors like Rubini et al. [XI] have suggested a variety of categorization systems. Using information on the correlation between diabetes and heart disease, the Random Forest Algorithm was used to determine the percentage of heart disease that might be predicted. An accuracy of 84.81 was provided by Random Forest. SVM, which has an accuracy rate of 64.4%, was chosen as the most effective method. Amin Ul Haq et al. [XII] implemented various existing ML algorithms probably 7 models. Among all models logistic regression gave the result and the accuracy is 89%.

Author Syed Nawaz Pasha et al. [XIII] used different deep learning algorithms like SVM, Decision Tree, KNN and ANN etc. for predicting heart disease. ANN turned out to be the best algorithm with an accuracy of 85.24%. Abhijeet Jagtap et al. [XIV] In their research, they used data mining techniques including SVM, Nave Bayes, and Logistic Regression to predict heart disease. Rohit Bharti et al. [XV] heart illness was predicted using a comparison of four machine learning methods. The authors concluded that when data preprocessing was applied, the KNN classifier performed better while using the Machine Learning technique. Bhavesh Dhande et al. [XVI] in their work predicted diabetes and heart disease using Machine Algorithms. The authors concluded that Sigmoid SVC gave the best results with an accuracy of 88%.

**Table 1: Similar Works in this field**

| Reference | Author | Sample size/ Time period | Technique | Findings/Conclusion |
|---|---|---|---|---|
| [V] | Patel, J., Tejal Upadhyay, D., and Patel, S. | Sep 2015 to March 2016 | J48, Logistic Model Tree Algorithm, Random Forest Algorithm | The J48 tree approach proved to be the most accurate classifier for predicting heart disease. |

*H. Singh et al.*

| [VI] | Soni, J., Ansari, Soni, S. U., Sharma, D., | June 2011 | Weighted Associative Classifier (WAC) | Weighted associative classifier-based approach used for heart disease prediction |
|---|---|---|---|---|
| [VII] | Y. E. Shao, C.-D. Hou, and C.-C. Chiu | March 2019 | Random forest algorithm, Naïve Bayes algorithm | With an accuracy of 86.81% for Dataset-1, an accuracy of 82.75% for Dataset-2, and an accuracy of 86.81% for Dataset-3, the Random Forest method exceeds etc. |
| [VIII] | V. Chauraisa and S. Pal | January 2014 | Support Vector Machine | The best classifier for CVD prediction, with a precision of 85.77%, was the support vector machine model |
| [IX] | Mrs.G. Subba lakshmi | April-May 2011 | Naïve Bayes | The model that best predicted heart disease patients was naive Bayes. |
| [X] | Sonam Nikhar, A.M. Karandikar | June- 2016 | Decision Tree and Naïve Bayes Classifier | For the prediction of cardiovascular diseases, many classification algorithms are applied. Compared with naive Bayes classifier, the decision tree was more accurate. |
| [XI] | Rubini PE, Dr.C.A. Subasini, V. Kumaresan, S. Gowdham Kumar, T.M. Nithya Dr.A. Vanitha Katharine, | January 2021-Feb 2021 | Random Forest, Linear Regression, SVM, Naïve Bayes | Using data on the correlation between diabetes and heart disease, the Random Forest model was applied to determine the percentage of the heart disease that might be predicted. Accuracy of 84.81 was provided by Random Forest. |
| [XII] | Amin Ulf Haqq, Jian Ping Li, Muhammad Hammad, Shah Nazir, and Ruin a Sun | Sep 2018- dec 2018 | KNN, logistic Regression, ANN, Naïve Bayes, SVM | The highest accuracy (89%) came from logistic regression. |
| [XIII] | Syed Nawaz Pasha, Dadi Ramesh, | 2020 | SVM, KNN, Decision Tree, ANN | Best Accuracy was given by ANN with 85.24%. |

*H. Singh et al.*

| | Sallauddin Mohmmad, A. Harshavardhan and Shabana | | | |
|---|---|---|---|---|
| [XIV] | Abhijeet Jagtap, Priya Malewadkar, Omkar Baswat, Harshali Rambade | Feb 2019 | SVM, Naïve Bayes, Logistic Regression | The most effective method, SVM, was chosen for its 64.4% accuracy. |
| [XV] | Rohit Bharti, Aditya Khamparia, Mohammad Shabaz, Gaurav Dhiman, Sagar Pande, and Parneet Singh | May 2021 -July 2021 | KNN,SVM, Decision Tree,Random Forest | When data preprocessing was used in the ML technique, KNeighbors classifier performed better. |
| [XVI] | Bhavesh Dhande, Kartik Bamble, Sahil Chavan, Tabassum Maktum | 2022 | Random Forest, KNN, Logistic Regression, Decision Tree, AdaBoost, SVC, XGBoost | Sigmoid SVC accuracy was best with 88%. |

## III.  Objectives

From the performance point of view, our model is compared with another ML algorithm to evaluate how well different machine-learning approaches, predict heart-related disorders.

## Dataset

This manuscript retrieved the data set from https://www.kaggle.com and it consists of 303 samples. The records include ["age", "sex", "cp", "trestbps", "chol", "fbs", "restecg", "oldpeak", "slope", "ca", "thal", and "target"], which represent various risk factors, demographic information, and medical history of the patients. The dataset was utilized to train and estimate the effectualness of ML algorithms for forecasting cardiac diseases. Below figure 1 shows the template of our dataset.

*H. Singh et al.*

| | 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 |
| 4 | 57 | 0 | 0 | 120 | 354 | 0 | 1 | 163 | 1 | 0.6 | 2 | 0 | 2 | 1 |

**Fig. 1.** Sample entries of the data set*.*

Figure 1 represents the sample data used in this work. It represents the first five entries of the dataset with their attributes.

The characteristics that were used to examine the results are shown in Figure 2. The Dataset used in this paper has a total of 303 rows and 14 data columns (13 features and one target value) to get better results.

```
 #    Column          Non-Null Count   Dtype
---   ------          --------------   -----
 0    age             303 non-null     int64
 1    sex             303 non-null     int64
 2    cp              303 non-null     int64
 3    trestbps        303 non-null     int64
 4    chol            303 non-null     int64
 5    fbs             303 non-null     int64
 6    restecg         303 non-null     int64
 7    thalach         303 non-null     int64
 8    exang           303 non-null     int64
 9    oldpeak         303 non-null     float64
 10   slope           303 non-null     int64
 11   ca              303 non-null     int64
 12   thal            303 non-null     int64
 13   target          303 non-null     int64
```

**Fig. 2.** Data Attributes

Figure 2 represents the data attributes where the respective datatype except the old peak all other datatypes are integer type whereas the other is of float type. The correlation between the attributes is given in Figure 3.
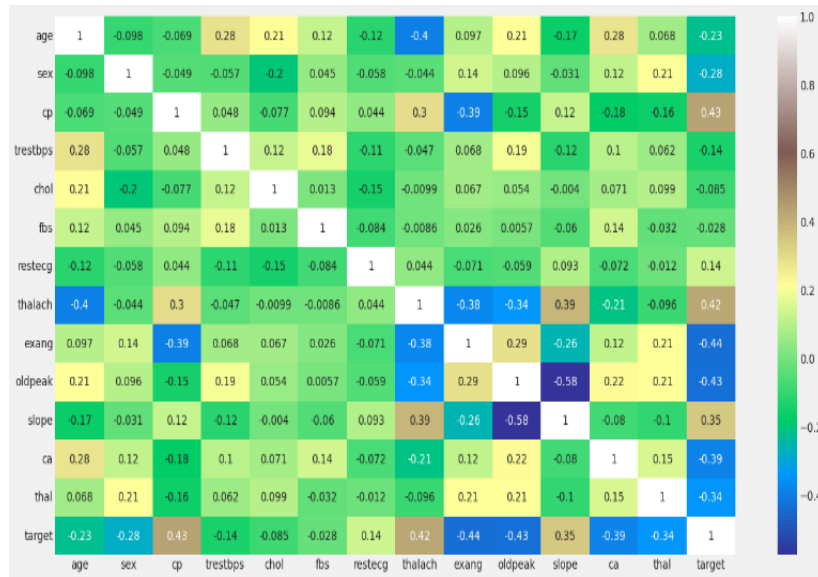
*H. Singh et al.*

**Fig. 3.** Correlation Matrix

Few characteristics have a negative correlation with the target value in the above correlation matrix, and even fewer have a positive correlation.
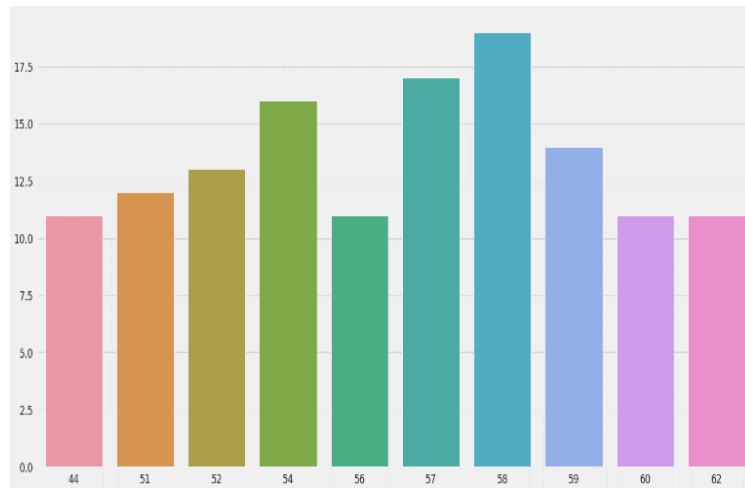The Pair plot of the attributes used is shown in Figure 4.



**Fig. 4.** Pair plot of the attributes used

The graph displayed above is called a Pair plot. Pair plots are used to find the most identifiable clusters or the most suitable set of traits to characterize a connection between two variables. It shows the relationship between one feature on the x-axis with all other features on the y-axis is shown in Figure 5.
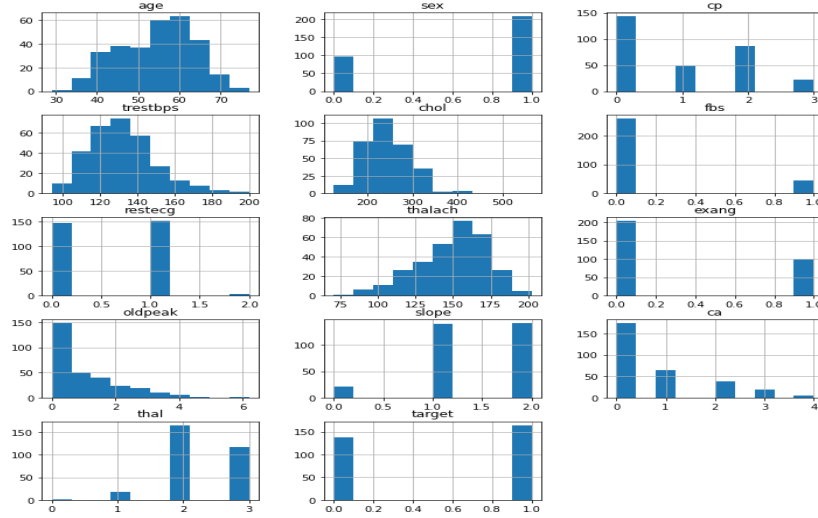
*H. Singh et al.*

**Fig. 5.** Histograms of all the attributes

The histograms for each variable are displayed above. Each characteristic has a distinct spectrum of distribution, as can be observed.
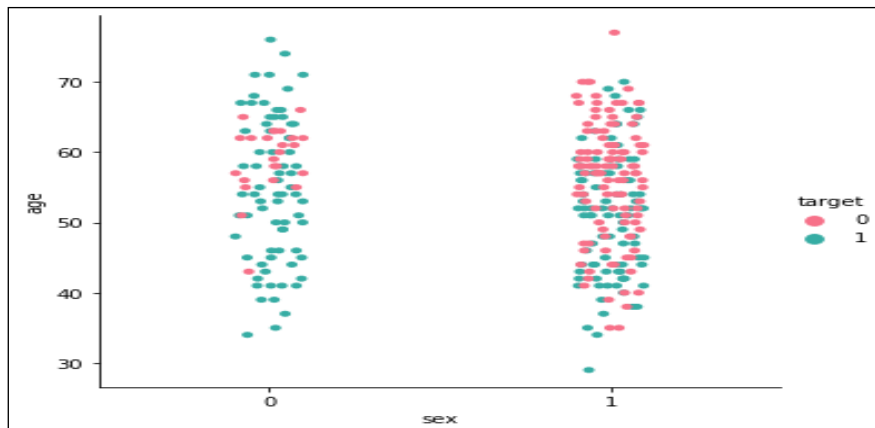


**Fig. 6.** Relation between age and sex.

The graph above shows the relationship between age and sex in terms of the presence of heart disease Figure 6. The oldest patient with no heart disease is a 77-year-old male and the youngest patient with a heart disease is a 29-year-old male patient
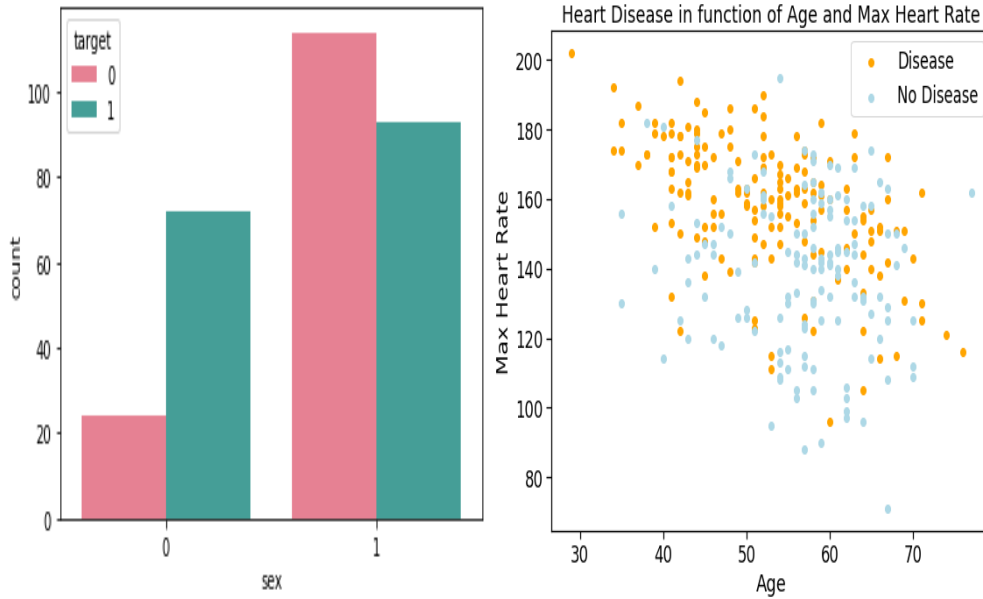
**Fig. 7.** Count of male/female with heart disease.  **Fig. 8.** Age and maximum heart rate in a scatter plot.

Here it is clearly visible that the slope value is higher in the case of males Figure 7. The number of healthy Male Patients is a lot higher than healthy Female Patients. Meanwhile, the number of Male Patients with heart disease is not that higher than Female Patients. The heart disease function of age and max heart rate is given in Figure 8. The figure is plotted by taking age as the x-axis and max-heart rate as the y-axis. However, from the figure it is observed that the max of heart disease exists at the rate of 200 age beyond 70.

## IV.  Methodology

This section describes various machine learning models implemented in this work. The methodology involves steps such as data preprocessing, model selection, model implementation, and model evaluation. The detailed methodology diagram is given in Figure 9.
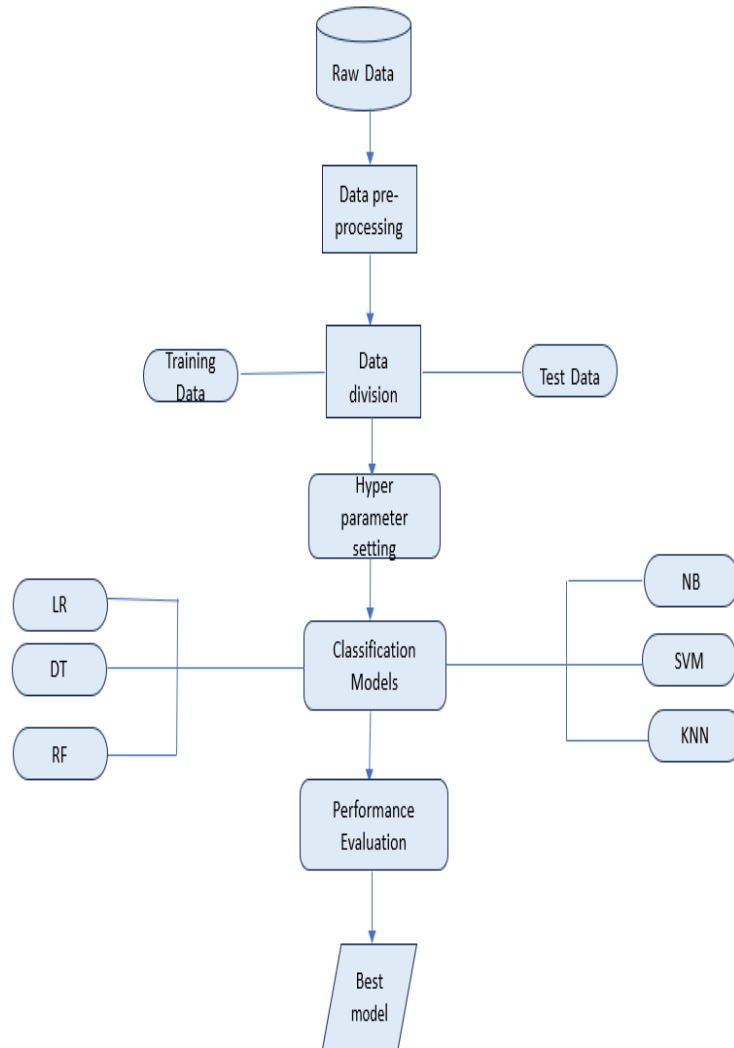
**Fig. 9.** Methodology diagram

## V. Implementations and Results Analysis

This section describes the implementation of six machine learning models and also analyzes the best model fit for the data selected for these experiments. The results of each experiment are explained through a confusion matrix followed by a classification report. The hyperparameter settings have been done to achieve the highest accuracy for each model.

**Logistic Regression-** The technique for binary classification issues is logistic regression. To predict a binary result, if the patient has heart disease or not, a linear model is used. The logistic function maps the input features to a probability value between 0 and 1, which can then be threshold to produce a binary prediction. To assess the importance of each factor in the prediction, the model's coefficients may be used,

to make informed decisions about which features are most relevant for the classification task. Additionally, Logistic Regression is computationally efficient and can handle high-dimensional data. The confusion matrix and classification report are given in Figure 10.
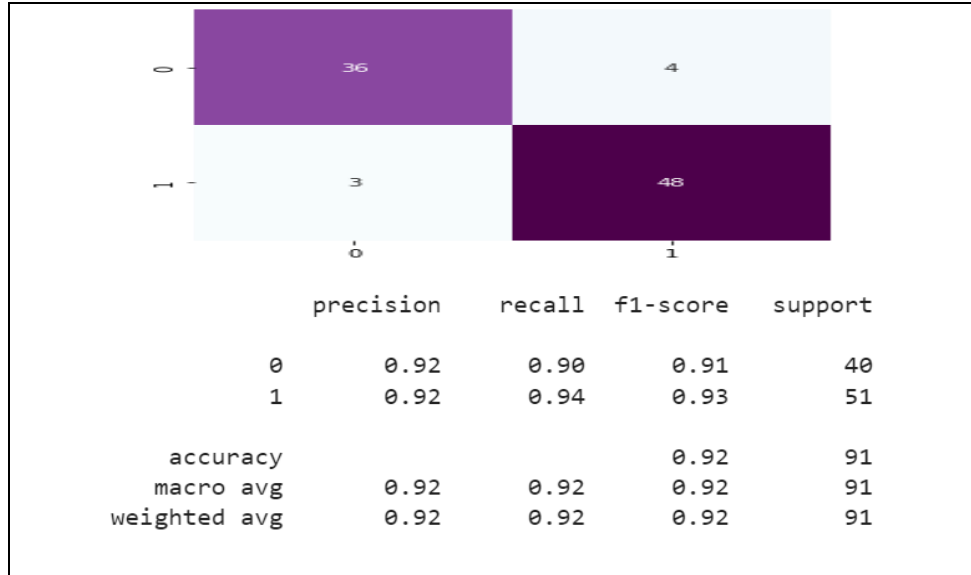


|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.92 | 0.90 | 0.91 | 40 |
| 1 | 0.92 | 0.94 | 0.93 | 51 |
| accuracy |  |  | 0.92 | 91 |
| macro avg | 0.92 | 0.92 | 0.92 | 91 |
| weighted avg | 0.92 | 0.92 | 0.92 | 91 |

**Fig. 10.** Confusion matrix and Classification report (LR model)

**Decision Tree –** The Decision tree model works on the principle of the structure of a tree, leaves, and nodes. Decision tree classifiers are very suitable for both classification and regression problems. The confusion matrix and classification report are given in Figure 11.
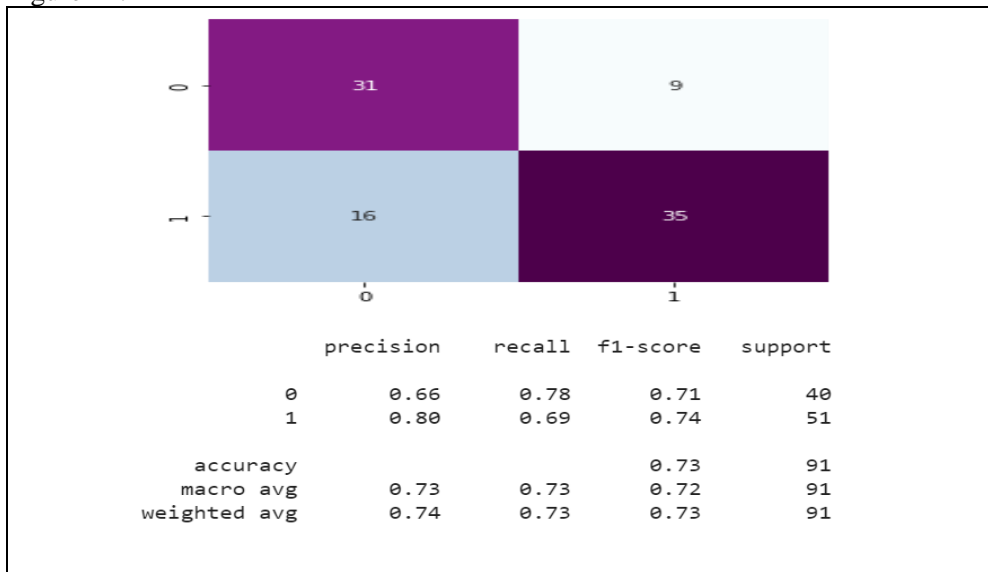


|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.66 | 0.78 | 0.71 | 40 |
| 1 | 0.80 | 0.69 | 0.74 | 51 |
| accuracy |  |  | 0.73 | 91 |
| macro avg | 0.73 | 0.73 | 0.72 | 91 |
| weighted avg | 0.74 | 0.73 | 0.73 | 91 |

**Fig. 11.** Confusion matrix and Classification report (DT model)

*H. Singh et al.*

**Random Forest -** It is a very popular AI technology that combines the results of multiple decision trees to reach an accurate result. Its purpose is to do tasks such as regression and classification. The confusion matrix and classification report are given in Figure 12.
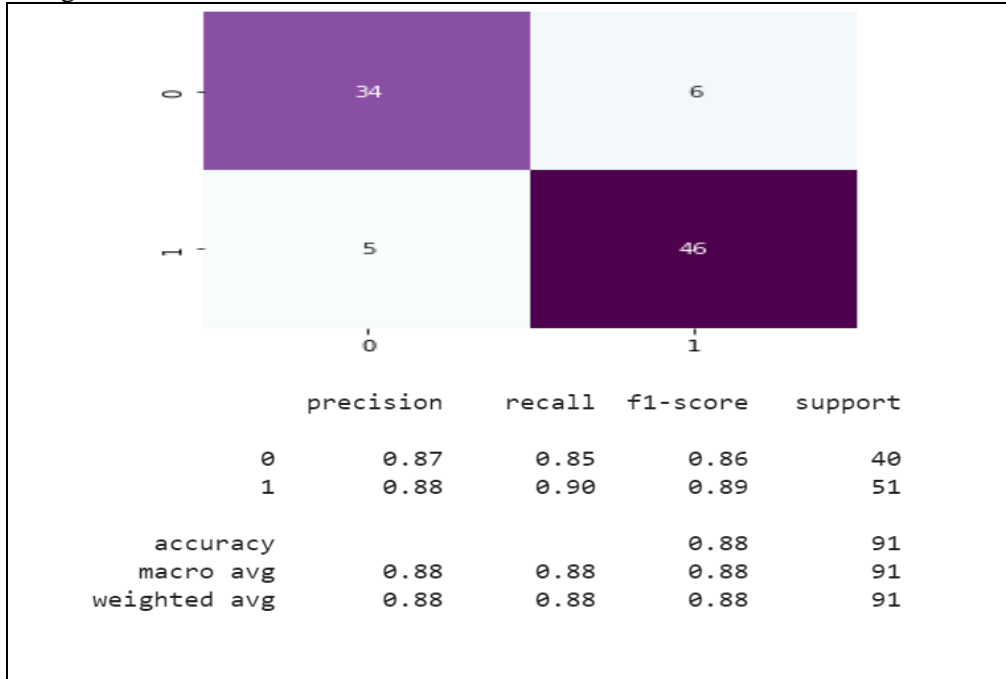


**Fig. 12.** Confusion matrix and Classification report (RF model)

Naïve Bayes - For classification issues, Naive Bayes, a probabilistic machine learning method, is used. The Bayes theorem, upon which this process is founded, asserts that the likelihood of an occurrence given the evidence is equal to the prior probability of the event occurring, multiplied by the prior probability of the event. In the context of heart disease prediction, the Naive Bayes algorithm assumes that the features (risk factors, demographic information, etc.) are independent of each other and calculates the probability of each class (heart disease or no heart disease) based on the evidence (features) of each patient. Naive Bayes has several advantages, one of which is its simplicity and computational efficiency. The algorithm is fast and easy to implement, making it a popular choice for many classification problems, including heart disease prediction.

$$\textbf{Equation – Bayes Theorem: } P(A/B) = \frac{P(A/B)P(A)}{P(B)} \tag{1}$$

In the case of the Bayes theorem, it is observed that both traits with predictors are independent. Thus, the presence of one trait does not affect the behaviour of another. The confusion matrix and classification report are given in Figure 13.
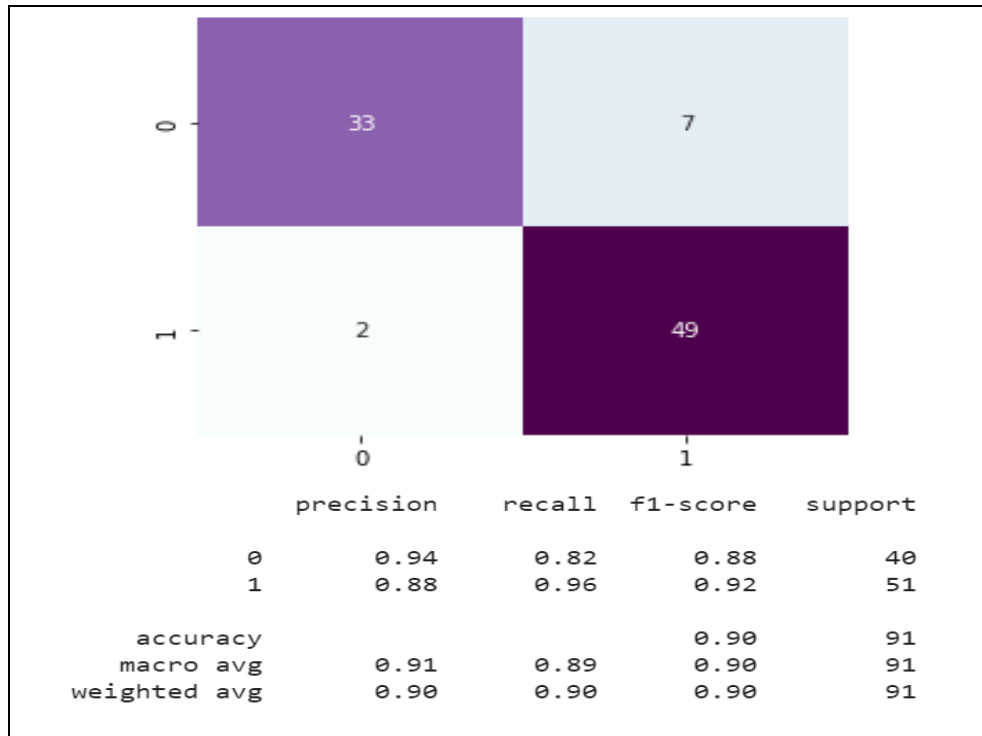
*H. Singh et al.*

**Fig. 13.** Confusion matrix and Classification report (NB model)

**Support Vector Machine (SVM)** – Nowadays for the solution of regression and classification problems SVM algorithm is widely used. It is a powerful tool for solving complex prediction problems and is particularly well-suited for problems where the data is not linearly separable. The hyperplane is selected in a manner that enhances the margin, which is the distance for each class between the nearest data points and the hyperplane. For heart disease prediction, SVM plays a vital role in determining whether a particular patient has cardiac disease or not, on diverse risk factors like demographic information, and medical history. The algorithm can be trained on a dataset of patients with known outcomes, and can then be used to predict the outcome for new patients. The confusion matrix and classification report are given in Figure 14.
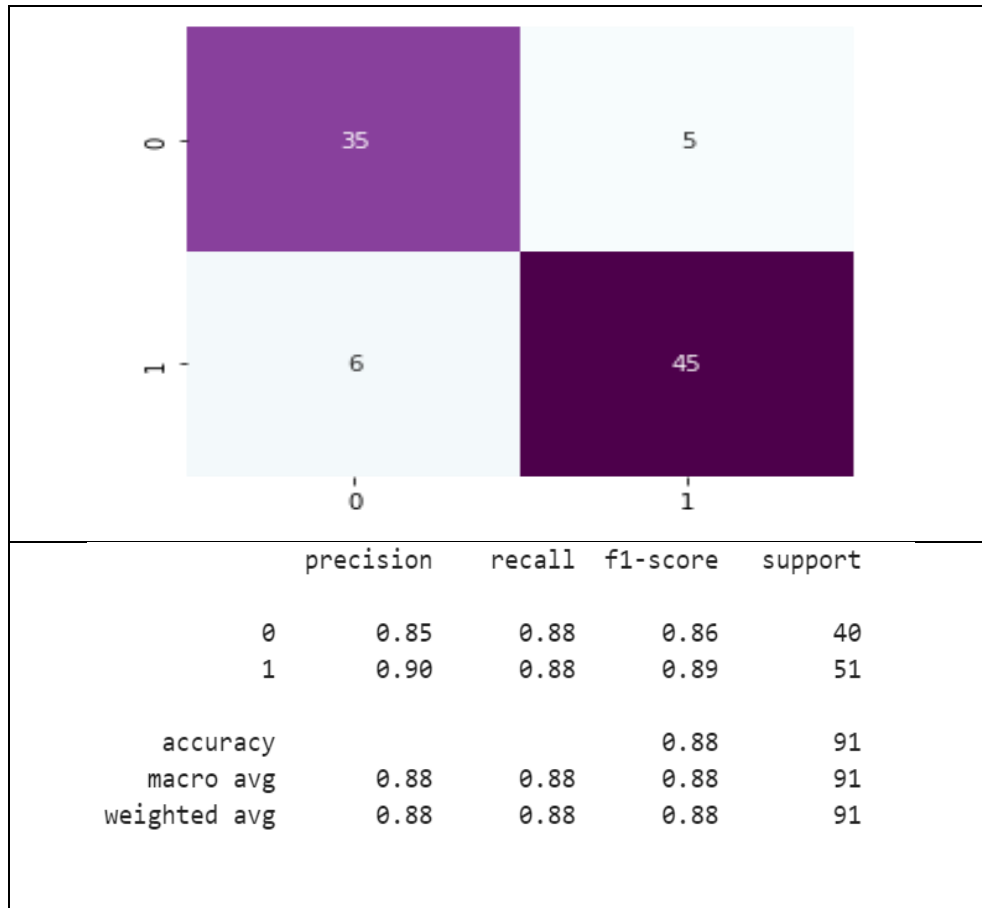
|   | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.85 | 0.88 | 0.86 | 40 |
| 1 | 0.90 | 0.88 | 0.89 | 51 |
| accuracy |  |  | 0.88 | 91 |
| macro avg | 0.88 | 0.88 | 0.88 | 91 |
| weighted avg | 0.88 | 0.88 | 0.88 | 91 |

**Fig. 14.** Confusion matrix and Classification report (SVM model)

**K-Nearest Neighbors (KNN)** - KNN is treated as a non-parametric ML technique employed for classification with regression applications. The majority class of an object's K nearest neighbours in the training dataset is used as the basis for classifying it according to the basic aspect of KNN. The ease with which KNN can handle non-linearly separable data is only one of its many benefits. The algorithm is also flexible, allowing for easy incorporation of domain knowledge and the ability to handle missing values and noisy data. One of the main challenges of KNN is choosing the value of K, this might have a substantial effect on the model's performance. A K number that is either too little or too big may lead to overfitting or underfitting, respectively. To find the ideal value of K, cross-validation is a popular method.

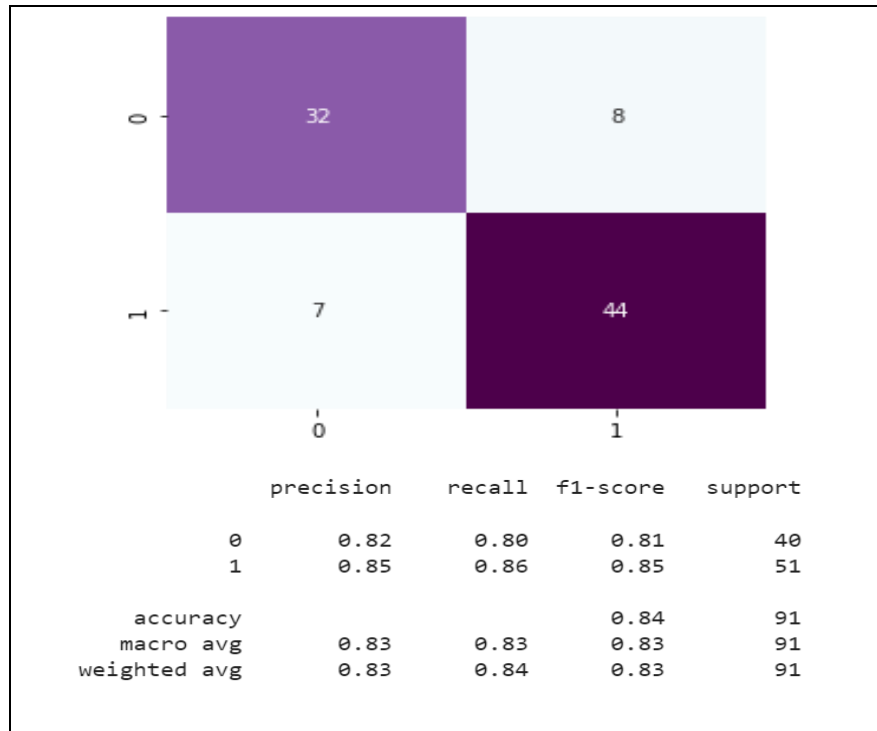The confusion matrix and classification report are given in Figure 15.

**Fig. 15.** Confusion matrix and Classification report (KNN model)

**Table 2: Best Accuracy for each Algorithm.**

| Algorithms | Accuracy |
|---|---|
| Logistic Regression | 92% |
| DT | 73% |
| RF | 88% |
| Naïve Bayes | 90% |
| SVM | 88% |
| KNN | 84% |

From the table, it can be observed that to predict heart disease, several algorithms are applied. The most accurate approach, according to the analysis of all the findings, is Logistic Regression with an accuracy of 92%.

## VI. Conclusion

In this work, four Machine learning algorithms are used to diagnose many types of heart diseases. From the experiment, the model found that Logistic Regression produces the best accuracy of 92% using our data set. However, the other algorithms

*H. Singh et al.*

such as DT, Random Forest, KNN, Naïve Bayes, SVM, and KNN have accuracy like 73%, 88%, 90%, 88%, 84% etc. 97 percent. The process of diagnosing CVD in the medical industry is expensive and time-consuming. The suggested method implies that Machine Learning may be utilized as a clinical tool in the identification of CVDs and will be especially helpful for doctors in the case of a false-positive result. In comparison to the previous methods described, the developed Machine learning models consistently predict illnesses with higher It is expected that the proposed method would enhance the medical sector. The suggested approach may be applied to categorize additional chronic illnesses, including thyroid, liver, breast, and diabetes mellitus. Using IoT and cloud computing approaches, the created algorithms may be used to huge data sets to predict such chronic illnesses. Based on the aforementioned study, it is clear that using ML approaches will significantly help in avoiding deaths and support medical professionals' efforts to reduce the onset of CVD among all patients. If adopted, this would be a prime example of how modern technology may be used for the good of everybody. Further research can be enhanced by identifying the most relevant elements needed for the prediction of heart disease and using other data mining techniques.

**Informed Consent:**

All authors have read and agreed to the published version of the manuscript and give their consent for the publication.

**References**

I. Kaur, B., Kaur, G., Heart Disease Prediction Using Modified Machine Learning Algorithm. International Conference on Innovative Computing and Communications. Lecture Notes in Networks and Systems, vol 473. Springer, Singapore. (2023)

II. V. Rawat, K. Gulati, U. Kaur, J.K Seth, V. Solanki, A.N. Venkatesh, D.P. Singh, N. Singh, M. Loganathan, "A Supervised Learning Identification System for Prognosis of Breast Cancer", Mathematical Problems in Engineering, vol. 2022, Article ID 7459455, 8 pages, 2022.

III. WHO Data - https://www.who.int/health-topics/cardiovascular-diseases#tab=tab_1

IV. Heart Disease Dataset - https://www.kaggle.com/c/heart-disease dated: Sept 2018

V. J. Patel, Tejal Upadhyay, D., and S. Patel, "Predicting heart condition using machine learning and data mining techniques". (2015)

*H. Singh et al.*

VI.   Soni, J., Ansari, U., Sharma, D., & Soni, S.  "Intelligent and effective heart disease prediction system using weighted associative classifiers". International Journal on Computer Science and Engineering. (2011)

VII.  Y. E. Shao, C.-D. Hou, and C.-C. Chiu, "Hybrid intelligent modelling, schemes for heart disease classification," Applied Soft Computing, vol. 14, pp. (2014).

VIII. V. Chauraisa and S. Pal, "Data Mining Approach to Detect  Heart Diseases,"  International Journal of Advanced Computer Science and Information Technology (IJACSIT), (2013).

IX.   Mrs. G. Subba lakshmi "Decision Support in Heart Disease Prediction System using Naive Bayes", Indian Journal of Computer Science and Engineering (IJCSE) (2011)

X.    Sonam Nikhar, and A. M. Karandikar. "Prediction of Heart Disease Using Machine Learning Algorithms." International Journal of Advanced Engineering, Management and Science, vol. 2, no. 6, Jun. 2016.

XI.   PE Rubini, Dr.C.A. Subasini, A. Vanitha Katharine, V. Kumaresan, S. Gowdham Kumar, T.M. Nithya. "A Cardiovascular Disease Prediction using Machine Learning Algorithms" Annals of R.S.C.B. (2021)

XII.  Amin Ul Haq, Jian Ping Li, Muhammad Hammad Memon, Shah Nazir, Ruinan Sun, "A Hybrid Intelligent System Framework for the Prediction of Heart Disease Using Machine Learning Algorithms", Mobile Information Systems, (2018).

XIII. Syed Nawaz Pasha, Dadi Ramesh, Sallauddin Mohmmad, A. Harshavardhan and Shabana "cardiovascular disease prediction using deep learning techniques" IOP Conf. Series: Materials Science and Engineering (2020)

XIV.  Abhijeet Jagtap, Priya Malewadkar, Omkar Baswat, Harshali Rambade "Heart Disease Prediction Using Machine Learning" International Journal of Research in Engineering, Science and Management (2019)

XV.   Rohit Bharti, Aditya Khamparia, Mohammad Shabaz, Gaurav Dhiman, Sagar Pande, Parneet Singh, "Prediction of Heart Disease Using a Combination of Machine Learning and Deep Learning", Computational Intelligence and Neuroscience (2021).

XVI.  Bhavesh Dhande, Kartik Bamble, Sahil Chavan, Tabassum Maktum "Diabetes & Heart Disease Prediction" ITM Web of Conferences ICACC-(2022)

*H. Singh et al.*