



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 10 Issue: VII Month of publication: July 2022

DOI: <https://doi.org/10.22214/ijraset.2022.44895>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Heart Disease Prediction Using Machine Learning Algorithms

Lakshmi. C. N¹, Bindhushree. M², Jaya Poojary³, Manish. C⁴, Shylaja. B⁵

BE Students, Asst. Professor, Department of Computer Science and Engineering, Dayananda Sagar Academy of Technology and Management, Bangalore, Karnataka, India

Abstract: Heart disease, alternatively known as cardiovascular disease, encases various conditions that impact the heart and is the primary basis of death worldwide over the span of the past few decades. It associates many risk factors in heart disease and a need of the time to get accurate, reliable, and sensible approaches to make an early diagnosis to achieve prompt management of the disease. Data mining is a commonly used technique for processing enormous data in the healthcare domain. Researchers apply several data mining and machine learning techniques to analyse huge complex medical data, helping healthcare professionals to predict heart disease. This research paper presents various attributes related to heart disease, and the model on basis of supervised learning algorithms as Naïve Bayes, decision tree, K-nearest neighbor, Support vector machine and random forest algorithm. It uses the existing dataset from the Cleveland database of UCI repository of heart disease patients. The dataset comprises 303 instances and 76 attributes. Of these 76 attributes, only 14 attributes are considered for testing, important to substantiate the performance of different algorithms. This research paper aims to envision the probability of developing heart disease in the patients.

I. INTRODUCTION

Over the last decade, heart disease or cardiovascular remains the primary basis of death worldwide. An estimate by the World Health Organization, that over 17.9 million deaths occur every year worldwide because of cardiovascular disease, and of these deaths, 80% are because of coronary artery disease and cerebral stroke [1]. The vast number of deaths is common amongst low and middle-income countries [2]. Many predisposing factors such as personal and professional habits and genetic predisposition accounts for heart disease. Various habitual risk factors such as smoking, overuse of alcohol and caffeine, stress, and physical inactivity along with other physiological factors like obesity, hypertension, high blood cholesterol, and pre-existing heart conditions are predisposing factors for heart disease. The efficient and accurate and early medical diagnosis of heart disease plays a crucial role in taking preventive measures to prevent death. Data mining is exploring huge datasets to extract hidden crucial decision making information from a collection of a past repository for future analysis. The medical field comprises tremendous data of patients. These data need mining by various machine learning algorithms. Healthcare professionals do analysis of these data to achieve effective diagnostic decision by healthcare professionals. Medical data mining using classification algorithms provides clinical aid through analysis. It tests the classification algorithms to predict heart disease in patients. Data mining is the process of extracting valuable data and information from huge databases. Various data mining techniques such as regression, clustering, association rule and classification techniques like Naïve Bayes, decision tree, random forest and K-nearest neighbor are used to classify various heart disease attributes in predicting heart disease. A comparative analysis of the classification techniques is used. In this research, I have taken dataset from the UCI repository. The classification model is developed using classification algorithms for prediction of heart disease. In this research, a discussion of algorithms used for heart disease prediction, comparison among the existing systems is made.

II. LITERATURE REVIEW

- 1) Rohit Bharti, Aditya Khamparia, Mohammad Shabaz, Gaurav Dhiman, Sagar Pande, and Parneet Singh "Prediction of Heart Disease Using a Combination of Machine Learning and Deep Learning" Volume 2021, Article ID 8387680, July 2021. Different machine learning algorithms and deep learning are applied to compare the results and analysis of the UCI Machine Learning Heart Disease dataset.
- 2) Devansh Shah, Samir Patel, Santosh Kumar Bharti "Heart Disease Prediction using Machine Learning Techniques", 16 October 2020. achieved 13 important attributes of HD that every researcher had done their research to predict the HD. we also achieved an appropriate two algorithm to provide high accuracy in heart disease dataset of the heart patient such as DT and NB and finally NB provides 88% accuracy among other algorithms. But the research of dataset of heart may failed due to distributed result of HD dataset.

- 3) Mangesh Limbitote, Pushkar Patil” A Survey on Prediction Techniques of Heart Disease using Machine Learning” ISSN: 2278-0181 Vol. 9 Issue 06, June-2020. In-depth analysis of relevant ml techniques to predict heart disease. Achieved accuracy of 82.30% using svm and 91.3% using Random forest.
- 4) Harshit Jindal,Sarthak Agrawal,Rishabh Khara “Heart disease prediction using machine learning algorithms.” , June 2020. Prepared heart disease prediction system to predict whether the patient is likely to be diagnosed with a heart disease or not using the medical history of the patient.
- 5) Dhai Eddine Salhi, Abdelkamel Tari, and M-Tahar Kechadi “Using machine learning for heart disease prediction” Feb 2021. Used data analytics to detect and predict disease’s patients. Starting with a preprocessing phase, where the most relevant features were selected then three data analytics techniques were applied on data sets of different sizes.
- 6) Baban.U. Rindhe, Nikita Ahire, Rupali Patil, Shweta Gagare, Manisha Darade “Heart Disease Prediction Using Machine Learning”, Volume 5, Issue 1, May 2021. The main objective of this research project is to predict the heart disease of a patient using machine learning algorithms

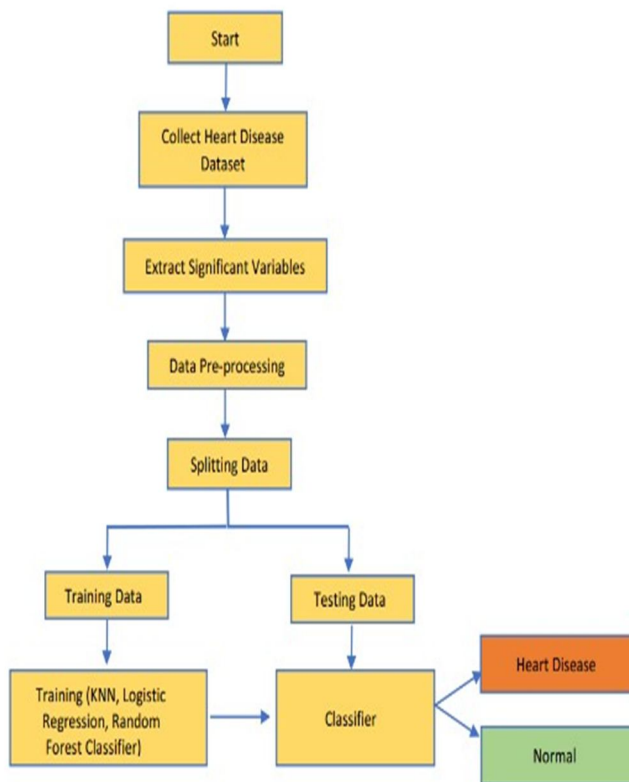
III. PROPOSED SYSTEM

A. Problem Statement

The problem statement is to develop a heart disease prediction system using machine learning. The system should extract knowledge and patterns which can help to get unbiased health decisions. The user should be able to input certain parameters of a patient based on which prediction will be made with high accuracy. The user interface should be minimal and easy so that it can be used by any non-technical person.

B. Methodology

Processed methodology start with the collection of data for this download tha data from kaggle that is well verified by researchers. In This methodology, There are many steps as shown in block diagram.



C. Description of the Dataset:

The dataset used for this research purpose was the Public Health Dataset and it is dating from 1988 and consists of four databases: Cleveland, Hungary, Switzerland, and Long Beach V. It contains 76 attributes, including the predicted attribute, but all published experiments refer to using a subset of 14 of them. The “target” field refers to the presence of heart disease in the patient. It is integer-valued 0 = no disease and 1 =disease. The first four rows and all the dataset features are shown in Table 1 without any preprocessing. Now the attributes which are used in this research purpose are described as follows and for what they are used or resemble:

- (i) Age—age of patient in years, sex—(1 =male; 0 =female).
- (ii) Cp—chest pain type.
- (iii) Trestbps—resting blood pressure (in mm Hg on admission to the hospital). The normal range is 120/80 (if you have a normal blood pressure reading, it is fine, but if it is a little higher than it should be, you should try to lower it. Make healthy changes to your lifestyle).
- (iv) Chol—serum cholesterol shows the amount of triglycerides present. Triglycerides are another lipid that can be measured in the blood. It should be less than 170 mg/dL (may differ in different Labs).
- (v) Fbs—fasting blood sugar larger than 120 mg/dl (1 true). Less than 100 mg/dL (5.6 mmol/L) is normal, and 100 to 125 mg/dL (5.6 to 6.9 mmol/L) is considered prediabetes.
- (vi) Restecg—resting electrocardiographic results.
- (vii) Thalch—maximum heart rate achieved. The maximum heart rate is 220 minus your age.
- (viii) Exang—exercise-induced angina (1 yes). Angina is a type of chest pain caused by reduced blood flow to the heart. Angina is a symptom of coronary artery disease.
- (ix) Oldpeak—ST depression induced by exercise relative to rest.
- (x) Slope—the slope of the peak exercise ST segment.
- (xi) Ca—number of major vessels (0–3) colored by fluoroscopy
- (xii) Thal—no explanation provided, but probably thalassemia (3 normal; 6 fixed defects; 7 reversible defects).
- (xiii) Target (T)—no disease = 0 and disease =1, (angiographic disease status).

Table 1 Attributes and details of dataset of heart disease

Sr. no.	Attribute	Representative icon	Details
1	Age	Age	Patients age, in years
2	Sex	Sex	0 = female; 1 = male
3	Chest pain	Cp	4 types of chest pain (1—typical angina; 2—atypical angina; 3—non-anginal pain; 4—asymptomatic)
4	Rest blood pressure	Trestbps	Resting systolic blood pressure (in mm Hg on admission to the hospital)
5	Serum cholesterol	Chol	Serum cholesterol in mg/dl
6	Fasting blood sugar	Fbs	Fasting blood sugar > 120 mg/dl (0—false; 1—true)
7	Rest electrocardiograph	Restecg	0—normal; 1—having ST-T wave abnormality; 2—left ventricular hypertrophy
8	MaxHeart rate	Thalch	Maximum heart rate achieved
9	Exercise-induced angina	Exang	Exercise-induced angina (0—no; 1—yes)
10	ST depression	Oldpeak	ST depression induced by exercise relative to rest
11	Slope	Slope	slope of the peak exercise ST segment (1—upsloping; 2—flat; 3—down sloping)
12	No. of vessels	Ca	No. of major vessels (0–3) colored by fluoroscopy
13	Thalassemia	Thal	Defect types; 3—normal; 6—fixed defect; 7—reversible defect
14	Num(class attribute)	Class	diagnosis of heart disease status (0—nil risk; 1—low risk; 2—potential risk; 3—high risk; 4—very high risk)

D. Data Pre-processing:

The real-life information contains large numbers with missing and noisy data. These data are pre-processed to overcome such issues and make predictions vigorously. Figure 1 explains the sequential chart of our proposed model. Cleaning the collected data usually has noise and missing values. To get an accurate and effective result, these data need to be cleaned in terms of noise and missing values are to be filled up. Transformation it changes the format of the data from one form to another to make it more comprehensible. It involves smoothing, normalization, and aggregation tasks.

Integration the data may not be acquired from a single source but varied sources, and it has to be integrated before processing. Reduction the data gained are complex and require to be formatted to achieve effective results. The data are then classified and split into training data set and test data set which is run on various algorithms to achieve accuracy score results.

IV. MODELS

A. Logistic Regression

Logistic regression is a classification algorithm used to assign observations to a discrete set of classes. Unlike linear regression which outputs continuous number values, logistic regression transforms its output using the logistic sigmoid function to return a probability value which can then be mapped to two or more discrete classes

```
In [9]: ##### Logistic Regression #####
X = df.iloc[:, :-1].values
y = df.iloc[:, -1].values

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 0)

from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression()
classifier.fit(X_train, y_train)

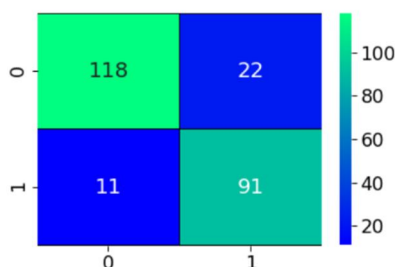
# Predicting the Test set results
y_pred = classifier.predict(X_test)

from sklearn.metrics import confusion_matrix
cm_test = confusion_matrix(y_pred, y_test)

y_pred_train = classifier.predict(X_train)
cm_train = confusion_matrix(y_pred_train, y_train)

print()
print('Accuracy for training set for Logistic Regression = {}'.format((cm_train[0][0] + cm_train[1][1])/len(y_train)))
print('Accuracy for test set for Logistic Regression = {}'.format((cm_test[0][0] + cm_test[1][1])/len(y_test)))

Accuracy for training set for Logistic Regression = 0.8636363636363636
Accuracy for test set for Logistic Regression = 0.892786885245902
```



B. Random forest [RF]

RF algorithm is supervised primarily based learning. It is used as classifier in numerous fields. By using this more trees makes a forest. If we have more number of trees then it create higher accuracy. It is also used for regression task. but it accomplish well when classify the task. And may overwhelmed misplaced values. There are three approach of RF: Forest RC(Random Blend) Forest RI(Random input) And combination of RC and RI.

Logistic regression [LR]:

LR is the supervised ML learning method. It is established on the association between dependent and independent variable as seen in Fig.5 variable “a” and “b” are dependent variable and independent variable and relation between them is shown by equation of line which is linear in nature that why this approach is called linear regression.

```
In [11]: ##### Random Forest #####
X = df.iloc[:, :-1].values
y = df.iloc[:, -1].values

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.4, random_state = 0)

from sklearn.ensemble import RandomForestClassifier
classifier = RandomForestClassifier(n_estimators = 10)
classifier.fit(X_train, y_train)

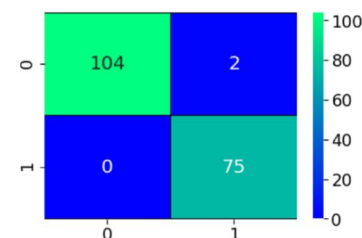
# Predicting the Test set results
y_pred = classifier.predict(X_test)

from sklearn.metrics import confusion_matrix
cm_test = confusion_matrix(y_pred, y_test)

y_pred_train = classifier.predict(X_train)
cm_train = confusion_matrix(y_pred_train, y_train)

print()
print('Accuracy for training set for Random Forest = {}'.format((cm_train[0][0] + cm_train[1][1])/len(y_train)))
print('Accuracy for test set for Random Forest = {}'.format((cm_test[0][0] + cm_test[1][1])/len(y_test)))

Accuracy for training set for Random Forest = 0.988950276243004
Accuracy for test set for Random Forest = 0.7459016393442623
```



C. Support Vector Machine

SVM is one type of ML method that work on the conception of hyper plan. It is used to find a hyper plan in n dimensional space, using this data point can be classified specifically [13]. (X_a, Y_a) is training sample of data set where $a=1,2,3,\dots,n$ and Y_a is the target vector and X_a is the i th vector. Hyper plan quantity select the variety of support vector such as example if a line is used as hyper plan then method is called linear support vector.

```
In [6]: ##### SVM #####
from sklearn.svm import SVC
classifier = SVC(kernel = 'rbf')
classifier.fit(X_train, y_train)

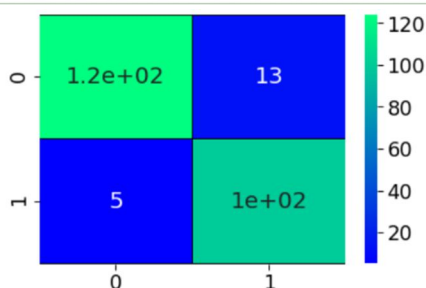
# Predicting the Test set results
y_pred = classifier.predict(X_test)

from sklearn.metrics import confusion_matrix
cm_test = confusion_matrix(y_pred, y_test)

y_pred_train = classifier.predict(X_train)
cm_train = confusion_matrix(y_pred_train, y_train)

print()
print('Accuracy for training set for svm = {}'.format((cm_train[0][0] + cm_train[1][1])/len(y_train)))
print('Accuracy for test set for svm = {}'.format((cm_test[0][0] + cm_test[1][1])/len(y_test)))

Accuracy for training set for svm = 0.9386792452830188
Accuracy for test set for svm = 0.8241758241758241
```



D. Decision Tree[DT]

DT is an algorithm that classifies parameters in categorical form in spite of arithmetic data. Tree like structure is created by DT. Many large data set related to medical have analyzed by DT due to its simple nature. It works on tree node for analysis. Leaf Node: Signify the solution of every Test Interior Node: Handle numerous element Main Node[Root Node]: Other nodes work based on main node Data is to be divided into two or more parallel set by applying this algorithm. Then entropy of each parameter is calculated. After that divide the data with predictor having extreme information gain that means minimum entropy

$$\text{Entropy}(S) = \sum_{i=1}^c -P_i \log_2 P_i,$$

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} |S_v| |S| \text{Entropy}(S_v).$$

```
In [10]: ##### Decision Tree #####
X = df.iloc[:, 1:].values
y = df.iloc[:, 0].values

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 0)

from sklearn.tree import DecisionTreeClassifier
classifier = DecisionTreeClassifier()
classifier.fit(X_train, y_train)

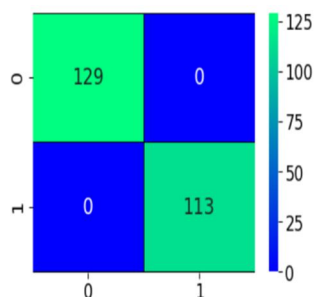
# Predicting the Test set results
y_pred = classifier.predict(X_test)

from sklearn.metrics import confusion_matrix
cm_test = confusion_matrix(y_pred, y_test)

y_pred_train = classifier.predict(X_train)
cm_train = confusion_matrix(y_pred_train, y_train)

print()
print('Accuracy for training set for Decision Tree = {}'.format((cm_train[0][0] + cm_train[1][1])/len(y_train)))
print('Accuracy for test set for Decision Tree = {}'.format((cm_test[0][0] + cm_test[1][1])/len(y_test)))

Accuracy for training set for Decision Tree = 1.0
Accuracy for test set for Decision Tree = 0.8932786885245902
```



E. Naïve bayes[NB]

NB is a supervise classification algorithm. It is a simple technique using Bayes theorem. To get the probability, mathematical concept is used with the support of bayes theorem. The correlation is neither related to each other nor predictor to one another. All parameters work autonomously for getting the maximum probability. $P(x/y) = P(Y/X) P(x) / p(y)$ Where $p(x)$ =Class predictor probability, $p(y)$ = Predictor Probability,

$P(x/y)$ = Posterior probability,

$P(y/x)$ =possibility, probability of predictor

```
In [7]: ##### Naive Bayes #####
X = df.iloc[:, :-1].values
y = df.iloc[:, -1].values

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 0)

from sklearn.naive_bayes import GaussianNB
classifier = GaussianNB()
classifier.fit(X_train, y_train)

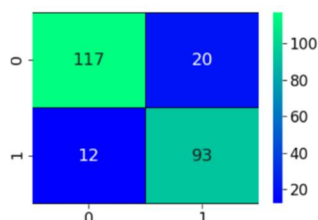
# Predicting the Test set results
y_pred = classifier.predict(X_test)

from sklearn.metrics import confusion_matrix
cm_test = confusion_matrix(y_pred, y_test)

y_pred_train = classifier.predict(X_train)
cm_train = confusion_matrix(y_pred_train, y_train)

print()
print('Accuracy for training set for Naive Bayes = {}'.format((cm_train[0][0] + cm_train[1][1])/len(y_train)))
print('Accuracy for test set for Naive Bayes = {}'.format((cm_test[0][0] + cm_test[1][1])/len(y_test)))

Accuracy for training set for Naive Bayes = 0.8677685950413223
Accuracy for test set for Naive Bayes = 0.7868852459016393
```



V. FLOWCHARTS

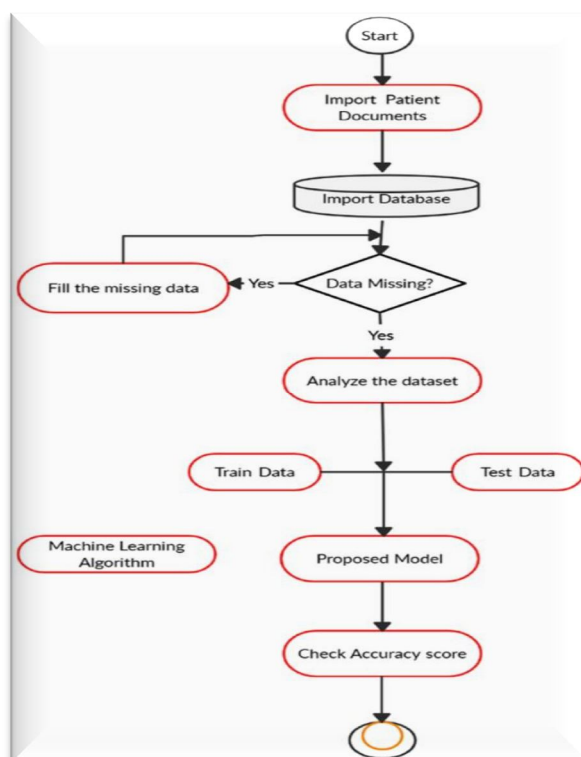


Figure.5.1. System Architecture

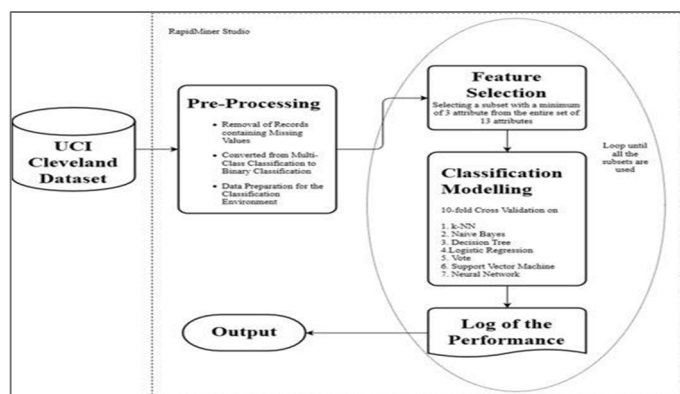


Figure 5.2- Dataflow process

VI. RESULT AND ANALYSIS

By applying different machine learning algorithms and then using deep learning to see what difference comes when it is applied to the data, three approaches were used. In the first approach, normal dataset which is acquired is directly used for classification, and in the second approach, the data with feature selection are taken care of and there is no outliers detection. The results which are achieved are quite promising and then in the third approach the dataset was normalized taking care of the outliers and feature selection; the results achieved are much better than the previous techniques, and when compared with other research accuracies, our results are quite promising.



Figure.6.3-R

MODULE	ACCURACY
Logistic Regression	86.3%
Random Forest Regression	98.8%
Support Vector Machine	93.8%
Decision Tree	100%
Naives Bayes	86.7%

VII. CONCLUSION

This project provides the deep insight into machine learning techniques for classification of heart diseases.

The role of classifier is crucial in healthcare industry so that the results can be used for predicting the treatment which can be provided to patients. The existing techniques are studied and compared for finding the efficient and accurate systems.

Machine learning techniques significantly improves accuracy of cardiovascular risk prediction through which patients can be identified during an early stage of disease and can be benefitted by preventive treatment.

It can be concluded that there is a huge scope for machine learning algorithms in predicting cardiovascular diseases or heart related diseases.

Aim of this research is to predict whether or not a patient will develop heart disease. This research was done on supervised machine learning classification techniques using Naïve Bayes, decision tree, random forest, support vector machine and K-nearest neighbor on UCI repository. Various experiments using different classifier algorithms were conducted through the WEKA tool. Research was performed on 8th generation Intel Core i7 having an 8750H processor up to 4.1 GHz CPU and 16 GB ram. Dataset was classified and split into a training set and a test set. Pre-processing of the data is done and supervised classification techniques such as Naïve Bayes, decision tree, K-nearest neighbor, and random forest are applied to get accuracy score. The accuracy score results of different classification techniques were noted using Python Programming for training and test data sets.

REFERENCES

- [1] Rohit Bharti, Aditya Khamparia, Mohammad Shabaz, Gaurav Dhiman, Sagar Pande, and Parneet Singh "Prediction of Heart Disease Using a Combination of Machine Learning and Deep Learning" Volume 2021, Article ID 8387680, July 2021.
- [2] Devansh Shah, Samir Patel, Santosh Kumar Bharti "Heart Disease Prediction using Machine Learning Techniques", 16 October 2020
- [3] Mangesh Limbote, Pushkar Patil "A Survey on Prediction Techniques of Heart Disease using Machine Learning" ISSN: 2278-0181 Vol. 9 Issue 06, June-2020.
- [4] Harshit Jindal, Sarthak Agrawal, Rishabh Khera "Heart disease prediction using machine learning algorithms." , June 2020.
- [5] Dhai Eddine Salhi, Abdelkamel Tari, and M-Tahar Kechadi "Using machine learning for heart disease prediction" Feb 2021.
- [6] Baban.U. Rindhe, Nikita Ahire, Rupali Patil, Shweta Gagare, Manisha Darade "Heart Disease Prediction Using Machine Learning", Volume 5, Issue 1, May 20
- [7] Chala Beyene, P. K., 2018. "Survey on Prediction and Analysis the Occurrence of Heart Disease Using Data Mining Techniques. International Journal of Pure and Applied Mathematics, 118
- [8] pp. 165-174. 8) Gupta Y, Ahmed R K A, Kautish S K; Application of data mining and knowledge management in special reference to medical informatics: a review; International Journal of Medical Laboratory Research 2017
- [9] Amandeep Kaur and Jyoti Arora .(2019). "Heart Diseases Prediction using Data Mining Techniques: A survey" International Journal of Advanced Research in Computer Science , IJARCS



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)